

# Structure Inducing Pre-Training

Matthew B. A. McDermott<sup>1,2</sup>, Brendan Yap<sup>1</sup>, Peter Szolovits<sup>1</sup>, Marinka Zitnik<sup>2,3,4,‡</sup>

<sup>1</sup>Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA, 02139, USA

<sup>2</sup>Department of Biomedical Informatics, Harvard Medical School, Boston, MA, 02115, USA

<sup>3</sup>Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA

<sup>4</sup>Harvard Data Science Initiative, Cambridge, MA 02138, USA

‡Corresponding author. Email: marinka@hms.harvard.edu

**Language model pre-training and derived methods are incredibly impactful in machine learning. However, there remains considerable uncertainty on exactly why pre-training helps improve performance for fine-tuning tasks. This is especially true when attempting to adapt language-model pre-training to domains outside of natural language. Here, we analyze this problem by exploring how existing pre-training methods impose relational structure in their induced per-sample latent spaces—*i.e.*, what constraints do pre-training methods impose on the distance or geometry between the pre-trained embeddings of two samples  $x_i$  and  $x_j$ . Through a comprehensive review of existing pre-training methods, we find that this question remains open. This is true despite theoretical analyses demonstrating the importance of understanding this form of induced structure. Based on this review, we introduce a descriptive framework for pre-training that allows for a granular, comprehensive understanding of how relational structure can be induced. We present a theoretical analysis of this framework from first principles and establish a connection between the relational inductive bias of pre-training and fine-tuning performance. We also show how to use the framework to define new pre-training methods. We build upon these findings with empirical studies on benchmarks spanning 3 data modalities and ten fine-tuning tasks. These experiments validate our theoretical analyses, inform the design of novel pre-training methods, and establish consistent improvements over a compelling suite of baseline methods.**

# Main

The pre-training (PT)/fine-tuning (FT) learning paradigm (also known as transfer learning) has had tremendous impact on natural language processing (NLP) and related domains [2, 35, 72]. In NLP or NLP-derived PT/FT, we are given a dataset  $\mathbf{X} \in \mathcal{X}^{N_{\text{PT}}}$  and attempt to pre-train an encoder  $f_{\theta} : \mathcal{X} \rightarrow \mathcal{Z}$  which maps our domain of interest  $\mathcal{X}$  into a latent space  $\mathcal{Z}$ :  $f_{\theta} : x_i \mapsto z_i$ . This encoder  $f_{\theta}$  is then transferred for use in various fine-tuning tasks (which are not known at pre-training time). We evaluate PT/FT systems via the transfer performance of  $f_{\theta}$  on said fine-tuning tasks.

In this work, we are concerned primarily with the efficacy of PT/FT for downstream tasks that operate at a *per-sample* level (e.g., in natural language processing, evaluating the sentiment of a whole restaurant review is a *per-sample* task, in contrast to identifying a named entity token within a sentence which is an *intra-sample/per-token* task). One aspect of pre-training that drives such eventual fine-tuning performance is the induced geometry of the pre-trained, per-sample latent space  $\mathcal{Z}$  (formally defined in the Methods section). For example, it is well documented that the sentence embeddings produced by pre-trained language models in NLP can be non-smooth and anisotropic, which harms downstream task performance [73]. In other domains, such as biomedical modalities, where per-sample tasks are even more prevalent than intra-sample tasks as compared to NLP, the importance of this geometry only increases. Despite this importance, research into mechanisms to induce explicit, deep structural constraints in  $\mathcal{Z}$  is surprisingly limited. Many methods outright ignore the geometry of  $\mathcal{Z}$  (e.g., by imposing no pre-training loss over the whole-sample embeddings during pre-training) [2, 4, 5, 5] and other methods impose either only shallow structural constraints, such as through an auxiliary, per-sample, classification PT objective [35, 40, 42], or deeper structural constraints, but in an implicit manner, such as through data-augmentation [56, 60] or noising-based contrastive losses [57, 59]. While such methods can be powerful and have been successful in many areas, we argue that the lack of a clear framework to design PT methods that impose structural constraints on  $\mathcal{Z}$  that are simultaneously *explicit* (similar to supervised classification losses) and *deep* (similar to noising/augmentation-based contrastive losses) is a major weakness.

On the basis of this observation, we develop an analytical framework under which the PT objective is subdivided into two components: first, a language-model inspired imputation/denoising objective that leverages intra-sample relationships, and, second, a loss term explicitly driven to regularize the geometry of the per-sample latent space  $\mathcal{Z}$  to reflect the connectivity patterns of a user-specified graph  $G_{\text{PT}}$ . By relying on graphs to capture the structure we wish to induce in  $\mathcal{Z}$ , this PT framework allows us to specify PT methods that induce *deep* structure in an *explicit* manner, filling exactly the gap identified above. In addition, this paradigm can capture diverse relationships, such as those motivated by external knowledge (e.g., [74]), self-supervised constraints (e.g., [75, 76]), or distances between samples in an alternate modality (e.g., [69]). Moreover, this

PT framework is simultaneously specific enough to allow us to make theoretical guarantees about how different PT graphs impact FT performance, general enough to encompass a variety of existing PT methods, and expressive enough to motivate new PT methods that have not been previously studied. In addition to theoretical analysis, we demonstrate empirically that defining new methods according to our framework, using explicit forms of real-world structure, yields significant benefits over competitive PT baselines across 3 modalities and 10 FT tasks.

Our work advances PT/FT research through three major contributions. First, we show via a comprehensive review and detailed commentary that existing pre-training methods largely do not induce structural constraints over  $\mathcal{Z}$  that are simultaneously *deep* and *explicit*. Second, we establish a new framework for describing PT methods, which provides a vehicle to design new PT methods that explicitly induce deep structural constraints in  $\mathcal{Z}$  in accordance with a user-specified PT graph  $G_{\text{PT}}$ . We further support this framework with theoretical results quantifying how the graph’s structure relates to FT task performance. Crucially, this formalization in our new PT paradigm offers insight into when PT does or does not add value over supervised learning alone. Third, we show that structure-inducing PT methods through our framework perform at or above the level of existing PT baselines across three data modalities and 10 FT tasks.

## Results

### General Pre-Training Problem Formulation

Given a dataset  $\mathbf{X}_{\text{PT}} \in \mathcal{X}^{N_{\text{PT}}}$ , a PT method aims to learn an encoder  $f_{\theta} : \mathcal{X} \rightarrow \mathcal{Z}$  such that  $f_{\theta}$  can be transferred to FT tasks that are unknown at pre-training time. While we can leverage additional information at PT time to inform the training of  $f_{\theta}$  (e.g., PT-specific labels  $\mathbf{Y}_{\text{PT}}$ ), the encoder  $f_{\theta}$  must take only samples from  $\mathcal{X}$  as inputs so that it can be used for fine-tuning. Pre-training methods typically solve this problem by training  $f_{\theta}$  to minimize a pre-training loss  $\mathcal{L}_{\text{PT}}$  over  $\mathbf{X}_{\text{PT}}$ . For example, in BERT,  $\mathcal{X}$  consists of free-text samples,  $f_{\theta}$  is a transformer model, and  $\mathcal{L}_{\text{PT}}$  consists of both a masked language modelling (MLM) per-token loss and the next-sentence-prediction (NSP) per-sample loss [35].

Note that our definition of pre-training ignores secondary applications of the pre-training objective itself; for example, autoregressive language models (e.g., GPT-3 [2]) are often used for their generative use directly, and not as commonly used to acquire embeddings or in transfer learning. This is a perfectly valid use of pre-trained language models within NLP, but is often not as useful in other domains which lack NLP’s generative properties, so we focus on the induced embeddings produced by pre-training methods instead. Note further that we are primarily interested in PT methods that either are or are derived from NLP PT methods. This domain is of particular interest because these methods (1) have been extremely successful within NLP [2, 35, 77], (2) have motivated a large number of derived methods in non-language, biomedical modalities [19, 33, 43, 46],

and (3) are not yet fully technically understood [29, 73, 78].

## Defining Explicit and Deep Structural Constraints

Central to our hypothesis is the claim that most NLP-derived PT methods today do not impose explicit, deep constraints on the (per-sample) latent space geometry of  $\mathcal{Z}$ . To justify this claim, we define “explicit” and “deep” structural constraints (Definitions 1-2).

### Definition 1. Explicit vs. Implicit Structural Constraints:

A PT objective  $\mathcal{L}_{PT}$  imposes a structural constraint that is *explicit* (vs. implicit) to the degree that it (as  $f_\theta$  approaches optimality) permits us to reason directly about the relationship (in particular, the distance) between any two samples  $z_i$  and  $z_j$  in the latent space  $\mathcal{Z}$ .

### Definition 2. Deep vs. Shallow Structural Constraints:

A PT objective  $\mathcal{L}_{PT}$  imposes a structural constraint that is *deep* (vs. shallow) on the basis of how much information (e.g., how many dimensions) would be required to fully satisfy the constraint.

For example, consider a classification PT loss according to labels  $y_i \in \mathcal{Y}$  and a logit layer which maps  $z_i \mapsto \tilde{y}_i$ . This method produces an *explicit* structural constraint because near optimality, we can infer that the relative (cosine) distance between two samples  $z_i$  and  $z_j$  is small if and only if  $y_i = y_j$ . However, this constraint is also *shallow*, because to fully satisfy this constraint, we need only embed each class  $c \in \mathcal{Y}$  with a unique position  $p_c \in \mathcal{Z}$ , then compress all samples  $z_i$  near their class prototype  $p_{y_i}$ . This distance-based constraint can be accomplished in a very low dimensional space  $\mathcal{Z}$  (e.g. we can distribute each  $p_c$  uniformly about a 2D unit circle, then compress all  $z_i$  to appear at a very small cosine distance from their class prototypes), illustrating that this constraint is very shallow.

In contrast, consider a contrastive method that asserts that  $z_i = f_\theta(x_i)$  should be close to  $z'_i = f_\theta(\tilde{x}_i)$ , under some noising/augmentation procedure  $x_i \mapsto \tilde{x}_i$ , but simultaneously far from other samples  $z_j$ . While this method constrains the latent space to be smooth with respect to the noising process, it offers only an *implicit* constraint on  $\mathcal{Z}$  as it is generally not possible to infer how the distance between distinct samples  $z_i$  and  $z_j$  is constrained. However, it imposes a *deeper* constraint than does the classification objective because the implicit connections between samples induced by the noising procedure reflect relationships that can not necessarily be captured in a low-dimensional space (dependent on dataset size and density).

## Existing Pre-training Methods do not use Deep, Explicit Constraints

To show that existing methods largely do not provide means to impose structural constraints that are simultaneously deep and explicit, we survey over 90 existing PT methods on the basis of how their objective functions constrain the  $\mathcal{Z}$  (Figure 1, Appendix A). For full details on our review

findings, see the Methods section. Throughout all examined methods, we find that *deep, explicit structural constraints are almost never employed*. Instead, most methods either (1) impose no per-sample PT objectives at all (*e.g.*, text-generation models, which are often not used for embeddings at all but rather for prompting or generative applications [2, 4–6]), (2) use explicit, but shallow, supervised PT objectives (*e.g.*, BERT’s “Next-sentence Prediction” (NSP) objective, ALBERT’s “Sentence-order Prediction” (SOP) objective, or various multi-task objectives [35, 40, 42]), or (3) use implicit, but deep, un- or self-supervised contrastive PT objectives (*e.g.*, contrastive sentence embedding losses [56, 57, 59, 60, 79]).

Across all surveyed methods, we find that only four methods impose simultaneously explicit and deep constraints: KEPLER [68], CK-GNN [69], XLM-K [70], and WebFormer [71]. All four can be described as some form of per-sample graph alignment, in which an external, pre-training knowledge graph  $G_{PT}$  or connectivity algorithm is employed over a subset of pre-training samples, and the output embeddings of pairs of samples  $z_i = f_{\theta}(x_i)$  and  $z_j = f_{\theta}(x_j)$  are constrained to reflect their relationships in the pre-training graph. This form of constraint is explicit, as the graph  $G_{PT}$  contains explicit relationships that will be induced in the output latent space, but also deep, as the geometry of the graph  $G_{PT}$  can be arbitrarily complex.

However, all these methods have major limitations. In KEPLER and XLM-K, the per-sample embeddings are only constrained to a restricted set of samples corresponding to entity descriptions from a knowledge graph. As such, there are no constraints implied on the general domain free-text samples in  $\mathcal{X}$  alone [68, 70]. In CK-GNN, the graph connectivity is derived from a cluster-restricted 1-nearest-neighbor graph in an alternate modality’s distance space, which may offer a limited higher-order structure, and unlike the NLP approaches, this method has no intra-sample (*e.g.* per-token) pre-training task [69]. Finally, in WebFormer, the graph used is inferred from the structure of the HyperText Markup Language (HTML) underlying web-pages, and relationships are only constrained at the per-sample level for limited structural relationships within the HTML. Further, WebFormer is a specialized model specifically for processing web content (text and HTML elements), so their approach can’t be directly generalized to other domains [71]. Moreover, these methods explore only the particular contexts of their individual models. They offer no general framework for how to realize this style of deep, explicit per-sample constraints in other contexts, nor do they explore any theory on how these constraints relate to performance for fine-tuning tasks [68–71].

Overall, our review of pre-training methods establishes unequivocally that pre-training methods capable of providing explicit, deep structural constraints are significantly under-explored. Across all the methods we reviewed, only four methods leverage constraints are explicit and deep, all of which have significant limitations, and there is no general consensus on how to constrain the  $\mathcal{Z}$  explicitly and deeply. These findings motivate our new framework, which offers insight into how to realize deep, explicit structural constraints in pre-training models across diverse contexts

and provides theoretical guidance on how structural constraints relate to fine-tuning performance.

## New Pre-training Framework: Structure-Inducing Pre-training (SIPT)

Our pre-training problem framework includes two small, but important, differences from the standard formulation (Figure 2).

First, we assume that we have as an additional input to the PT problem a graph  $G_{PT} = (V, E)$  where vertices denote pre-training samples within  $\mathbf{X}_{PT}$  (e.g.,  $\{\mathbf{x}_{PT} | \mathbf{x}_{PT} \in \mathbf{X}_{PT}\} \subseteq V$ ) and edges represent user-specified relationships. Importantly, while we take the graph  $G_{PT}$  an input to the PT problem, *we cannot use it as a direct input to  $f_\theta$* . Just like in traditional pre-training,  $f_\theta$  must take as input only samples from  $\mathcal{X}$ . *This is because otherwise, we can not apply  $f_\theta$  to the same, general class of FT tasks over domain  $\mathcal{X}$ .*

Second, we decompose the PT loss  $\mathcal{L}_{PT}$  into two components, weighted with hyperparameter  $0 \leq \lambda_{SI} \leq 1$ :

$$\mathcal{L}_{PT} = (1 - \lambda_{SI})\mathcal{L}_M + \lambda_{SI}\mathcal{L}_{SI}.$$

$\mathcal{L}_M$  is a traditional, intra-sample objective (e.g., a language model), and  $\mathcal{L}_{SI}$  is a new, structure-inducing objective designed to regularize the per-sample latent space geometry in accordance with the relationships (edges) in  $G_{PT}$ . Under our framework,  $\mathcal{L}_{SI}$  is only allowable for  $G_{PT}$ ,  $f_\theta$ , and  $\mathcal{Z}$  if it permits some stable optima at which point a radius nearest-neighbor connectivity algorithm under some distance function in  $\mathcal{Z}$  will recover  $G_{PT}$  (formal constraint is in the Methods section). Note that this constraint strikes a connection between our framework and the wealth of existing research focused on *graph representation learning* [80–85]. These techniques do indeed offer valuable insights into how to sample minibatches over graph-structured data and devise losses for graph embeddings; however, many methods for actually modelling graph-structured data, including deep attributed graph embeddings and graph convolutional neural networks, should not be seen as replacements for our techniques here as they are typically not adaptable to contexts in which the graph is not known at inference time, and so *they could not be used in our pre-training setting where  $f_\theta$  must take in only inputs from  $\mathcal{X}$  directly.*

As the new loss term added  $\mathcal{L}_{SI}$  is explicitly designed to *induce the structure of  $G_{PT}$  in  $\mathcal{Z}$* , we call methods trained under our framework *structure-inducing pre-training (SIPT) methods*. Many existing PT approaches can be re-realized as SIPT methods, including classification-based PT objectives like NSP or SOP, contrastive methods, or existing graph alignment methods (see Methods for full details).

## Theoretical Analyses

Under our framework, one can link the structure of the PT graph  $G_{PT}$  to eventual FT task performance. In particular, as a SIPT embedder  $f$  over graph  $G_{PT}$  approaches optimality under the loss  $\mathcal{L}_{SI}$ , it produces an embedding space such that nearest-neighbor performance for any downstream

task is lower bounded by the performance that could be obtained via a nearest neighbor algorithm over graph  $G_{\text{PT}}$  (Theorem 1). This fact directly connects the geometry of the graph  $G_{\text{PT}}$  with the eventual fine-tuning performance of a SIPT embedder  $f$ . Furthermore, it demonstrates the advantage of employing an explicit constraint rather than an implicit one; by controlling the structure of  $G_{\text{PT}}$ , users can directly choose to add different inductive biases to the PT process, in a manner which has a provable impact on the eventual suitability for downstream FT tasks.

**Theorem 1.** Let  $\mathbf{X}_{\text{PT}}$  be a PT dataset,  $G_{\text{PT}}$  be a PT graph, and let  $f_{\theta^*}$  be an encoder pre-trained under a PT objective permissible under our framing that realizes a  $\mathcal{L}_{\text{SI}}$  value no more than  $\ell^*$ . Then, under embedder  $f$ , the nearest-neighbor accuracy for a FT task  $y$  converges as dataset size increases to at least the local consistency (Definition 5) of  $y$  over  $G_{\text{PT}}$ .

We also establish two important corollaries of Theorem 1 that further illustrate the importance of choosing graphs  $G_{\text{PT}}$  which impose *deep* structural constraints (Corollaries 1-2).

**Corollary 1.** Let  $\mathbf{X}_{\text{PT}} \in \mathcal{X}^N$ , be a PT dataset with corresponding labels  $\mathbf{y} \in \mathcal{Y}_{\text{PT}}^N$ . Define  $G_{\text{PT}} = (\mathbf{X}_{\text{PT}}, E)$  such that  $(\mathbf{x}_i, \mathbf{x}_j) \in E$  if and only if  $y_i = y_j$ .

Then, the local consistency for a given FT task  $\mathbf{y}^{(\text{FT})}$  over  $G_{\text{PT}}$  (and thus by Theorem 1, the nearest-neighbor accuracy for any optimized SIPT embedder) is upper bounded by the probability that a sample  $x_i$ 's fine-tuning label  $y_i^{(\text{FT})}$  agrees with the majority class label for task  $\mathbf{y}^{(\text{FT})}$  over the clique consisting of all nodes with the same pre-training label  $y_i$  as  $x_i$ .

**Corollary 2.** Let  $\mathbf{X}_{\text{PT}}$  be a PT dataset that can be realized over a valid manifold  $\mathcal{M}$ . Assume  $\mathbf{X}_{\text{PT}}$  is sampled with full support over  $\mathcal{M}$ . Let  $G_{\text{PT}}(\mathbf{X}_{\text{PT}}, E)$  be an  $r$ -nearest-neighbor graph over  $\mathcal{M}$  (e.g.,  $(\mathbf{x}_i, \mathbf{x}_j) \in E$  if and only if the geodesic distance between the two points on  $\mathcal{M}$  is less than  $r$ :  $\mathcal{D}_{\mathcal{M}}(\mathbf{x}_i, \mathbf{x}_j) < r$ ). Let  $y^{(\text{FT})}$  be a FT classification task that is almost everywhere smooth on the manifold.

Then, as PT dataset size (and thus the size of  $G_{\text{PT}}$ ) tends to  $\infty$ , and  $r$  tends to zero, the local consistency of  $y^{(\text{FT})}$  over  $G_{\text{PT}}$  (and thus by Theorem 1 the nearest-neighbor accuracy of an SIPT embedder) will likewise tend to 1.

Informally, these corollaries establish that when a shallow structural constraint is used (e.g. a supervised classification objective), then the associated SIPT-equivalent model permits only minimal guarantees for FT performance, driven by the extent to which an FT task label is consistent within the classes under the supervised PT objective. In contrast, if a deep structural constraint is used, realized in Corollary 2 via  $G_{\text{PT}}$  being a nearest-neighbor graph over an arbitrary manifold  $\mathcal{M}$ , then a SIPT model permits a theoretical guarantee for FT performance that approaches unity as the pre-training dataset size grows for any FT task that is smooth over  $\mathcal{M}$ .

In sum, this theoretical analysis shows that we can directly connect the structure induced in  $\mathcal{Z}$  to downstream FT performance. As such, moving to new PT methods which leverage graphs  $G_{\text{PT}}$

with deeper structural constraints has the potential to markedly improve performance, as we will demonstrate on real-world datasets in our experiments. Complete proofs for all theoretical results and semi-synthetic experiments validating our theoretical findings in practice are in the Methods section.

## Real-world Experiments: Datasets and Tasks

We examine three data modalities for our experiments: PROTEINS, containing protein sequences; ABSTRACTS, containing free-text biomedical abstracts; and NETWORKS, containing sub-graphs of protein-protein interaction (PPI) networks.

In each data modality, we use different pre-training datasets and leverage different kinds of pre-training graphs  $G_{PT}$ , test on publicly available benchmarks for FT tasks, and compare our SIPT methods to compelling baselines spanning both per-sample and/or per-token methods (Tables 1-3). Further details on these aspects can also be found in the Methods Section.

## Real-world Experiments: $\mathcal{L}_{SI}$ and Training Procedures

As discussed in the definition of our framework, a SIPT method differs from a standard PT method by (1) the choice of graph  $G_{PT}$  (Table 1) and (2) the design of the new, structure-inducing loss  $\mathcal{L}_{SI}$ . To define  $\mathcal{L}_{SI}$  in our experiments, we leverage ideas from *structure-preserving metric learning* (SPML) [86–88]. SPML is a form of metric learning where positive relationships are defined by edges in a graph rather than a shared supervised label. We adapt two losses, a traditional contrastive loss [89] and a multi-similarity loss [90], from supervised metric learning to the graph-based, structure-preserving context of  $\mathcal{L}_{SI}$  terms in SIPT.

In addition to these losses, in the ABSTRACTS and PROTEINS domains, we use a warm-start procedure to initialize pre-training from existing language models rather than beginning from scratch. This saves significant computational time and allows for a powerful ablation study to isolate performance improvements to the introduction of our  $\mathcal{L}_{SI}$  term. Second, we perform extensive hyperparameter tuning studies on these two domains to identify appropriate values for  $\lambda_{SI}$ , and adapt those findings to the NETWORKS domain. Further details about the experimental setup, including formal statements of our contrastive and multi-similarity losses, are in the Methods section.

## Result 1: Incorporating $\mathcal{L}_{SI}$ performs comparably to or improves over all baselines across all 3 domains and 10 FT tasks

To analyze our experiments, we compute the relative reduction of error of the best performing SIPT model vs. the per-token or per-sample baselines across all FT tasks (Table 2). *We can see that in 10/15 cases, SIPT improves over existing methods, and in no case does it do worse than either baseline.* In some cases, the gains in performance are quite significant, with improvements of

approximately 17% (0.05 macro-F1 raw change) on AA, 6% on SRE (0.01 macro-F1 raw change), and 4% on RH (2% accuracy raw change). *SIPT models further establish a new SOTA on AA and RH and match SOTA on FL, ST, & PF.*

We see in Figure 3 how performance evolves over FT iterations for the NETWORKS dataset to determine if the improvements observed at the final converged values are present throughout training. We see that SIPT methods converge faster to better performance than both baselines. Raw results across all settings are presented in the Methods section (Tables 7-8).

## **Result 2: These performance gains are present across diverse modalities and pre-training graphs and outperform both per-sample and per-token baselines**

SIPT performance gains persist over all three data modalities and all different  $G_{PT}$  types we use here. This shows that explicitly regularizing the per-sample latent space geometry offers value across NLP, non-language sequences, and non-sequential domains, as well as while leveraging graphs including those defined by external knowledge, by self-supervised signals in the data directly, and by nearest-neighbor methods over multi-task label spaces. *Furthermore, note that these improvements exist not only in comparison to standard language modelling approaches but also against existing methods that impose per-sample PT objectives, including single and multi-task classification objectives.*

## **Result 3: Observed gains are uniquely attributable to the novel loss $\mathcal{L}_{SI}$**

As outlined in the Methods section, our experimental design permits us to determine how much of the observed gains in Table 2 are due to the novel loss component, as opposed to, for example, continued training, new PT data, or the batch selection procedures used in our method which also indirectly leverage the knowledge inherent in  $G_{PT}$ . Unsurprisingly, some gains are observed due to these other factors, and performance gains shrink when considering these ablation studies. However, even when comparing against the maximal performance baseline or ablation study overall, neither the direction of observed relationships nor the statistical significance of observed comparisons changes. *Therefore, we can conclusively state that the performance improvements observed here are uniquely attributable to the new, structure-inducing components introduced by our framework.* Full ablation study results can be found in the Methods section (Tables 7-8).

## **Discussion**

We show that despite the breadth of research into PT methods, methods for imposing *explicit* and *deep* structural constraints over the per-sample, pre-training latent space  $\mathcal{Z}$  are under-explored (Figure 1). Our theoretical and empirical analyses *show that this deficit matters in practice.* In particular, we define a new pre-training framework, *structure-inducing pre-training* (SIPT), under which the PT loss is subdivided into two components: one which is designed to capture intra-

sample (*e.g.* per-token) relationships and one which is designed to constrain the per-sample latent space to capture relationships between samples given by a user-specified pre-training graph  $G_{PT}$ . Under our framework, we show both theoretically and via experiments that the structure induced in  $\mathcal{Z}$  can be directly connected to eventual fine-tuning performance. Empirically, we show that novel SIPT methods leveraging a variety of pre-training graphs can consistently outperform compelling existing PT methods across three real-world domains.

Our work highlights several important directions for future research. For example, are there losses better suited than metric learning losses for pre-training graphs—*e.g.*, can we leverage the graph distance alongside the intra-batch distance to improve negative sampling strategies? In addition, can we produce theoretical results on convergence of pre-trained models? Can we advance the understanding of when and how pre-trained models converge to solutions that recover  $G_{PT}$ ? In a different direction, can pre-trained models reflect forms of structure beyond nearest neighbor relationships—*e.g.*, such as by leveraging higher-order topological considerations or by matching a distance function rather than a discrete graph? We anticipate that further analyses of these and other questions will lead to new pre-training methods and enable pre-training to be successful across diverse domains.

**Data availability.** Our synthetic datasets and pointers to all real-world datasets used (which are all publicly available) are available here: [https://github.com/mmcdermott/structure\\_inducing\\_pre-training](https://github.com/mmcdermott/structure_inducing_pre-training).

**Code availability.** All code for this project is available at [https://github.com/mmcdermott/structure\\_inducing\\_pre-training](https://github.com/mmcdermott/structure_inducing_pre-training).

**Acknowledgements.** MBAM was partly supported by a National Institutes of Health (NIH) grant LM013337 and a collaborative research agreement with IBM. BY was supported by a Massachusetts Institute of Technology (MIT) Undergraduate Research Opportunity fund. MZ gratefully acknowledges the support by the NSF under Nos. IIS-2030459 and IIS-2033384, US Air Force Contract No. FA8702-15-D-0001, and awards from Harvard Data Science Initiative, Amazon Research, Bayer Early Excellence in Science, AstraZeneca Research, and Roche Alliance with Distinguished Scientists. Any opinions, findings, conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the funders.

**Authors contribution.** MBAM and BY collated datasets, wrote modelling code, and ran experiments. MBAM compiled final results and completed the review of existing pre-training studies. MBAM, PS, and MZ conceived of the study and shaped the framing of the work. PS and MZ offered insight and guidance throughout the project. MBAM and MZ wrote the final manuscript, and MBAM, BY, PS, and MZ contributed edits to manuscript drafts.

**Competing interests.** The authors declare no competing interests.

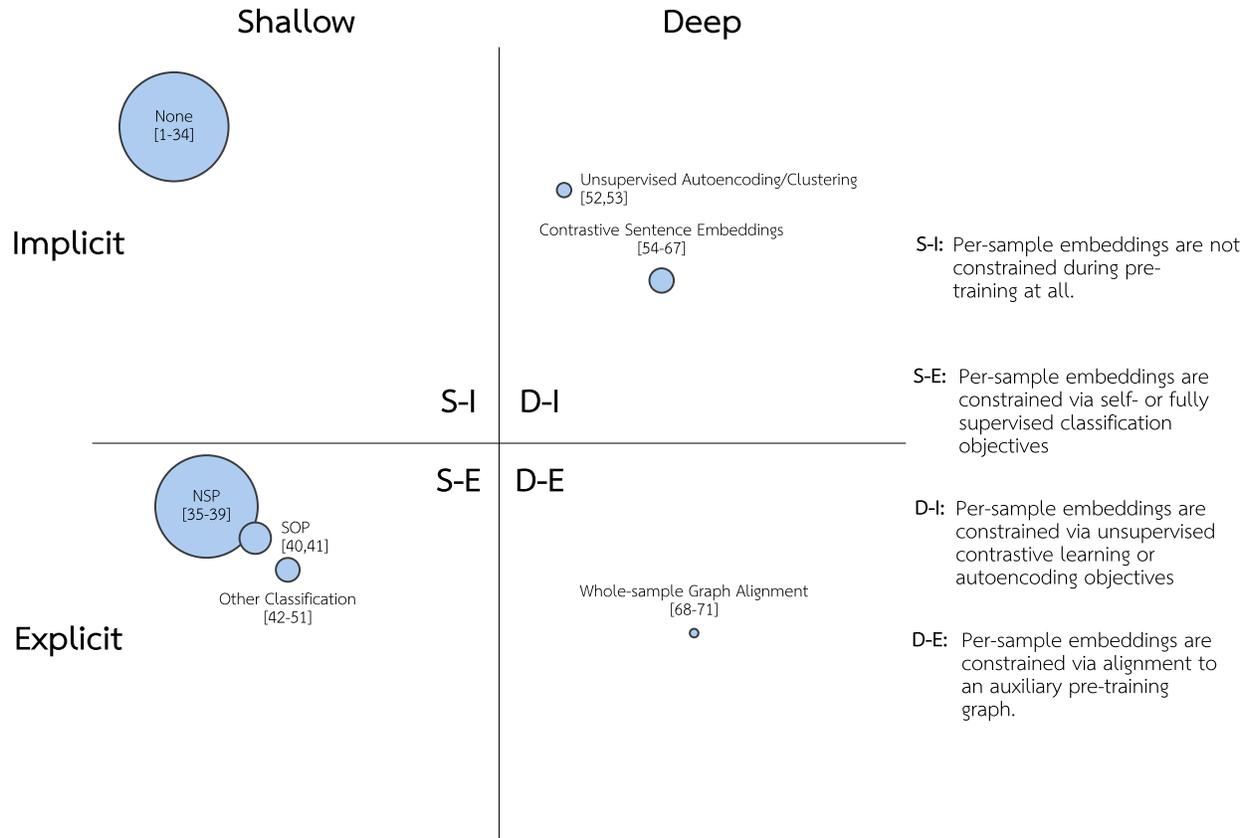


Figure 1: **Existing Pre-training (PT) Methods:** A summary of 71 existing natural language processing (NLP) and NLP-derived PT methods, categorized into clusters based on how they impose structural constraints over the PT (per-sample) latent space. Clusters are arranged on axes via manual judgements on whether the imposed constraint is *shallow* vs. *deep* and *implicit* vs. *explicit*. Clusters are sized such that the area corresponds to the number of citations methods included in that cluster have received on average per month since first publication, according to Google Scholar’s citation count. “None” captures models that leverage no pre-training loss over the per-sample embedding. “NSP” refers to “Next-sentence Prediction,” the per-sample PT task introduced in BERT [35]. “SOP” refers to “Sentence-order Prediction,” the per-sample PT task introduced in ALBERT [40]. Note that over 90 studies in total were considered in our review, but only 71 met the inclusion criteria to be included in this figure. These methods are described in more detail in Methods Table 4 and in Appendix A.

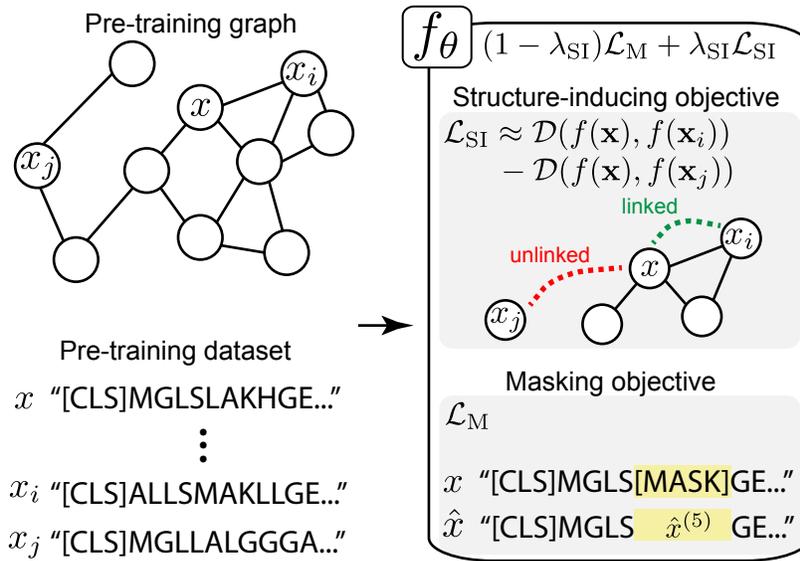


Figure 2: **Our Pre-training (PT) Framework:** We re-cast the PT formulation by taking a pre-training graph  $G_{PT}$  as an auxiliary input.  $G_{PT}$  is used to define a new structure-inducing objective  $\mathcal{L}_{SI}$ , which pushes a pre-training encoder  $f_{\theta}$  to embed samples such that samples are close in the latent space if and only if they are linked in  $G_{PT}$ .

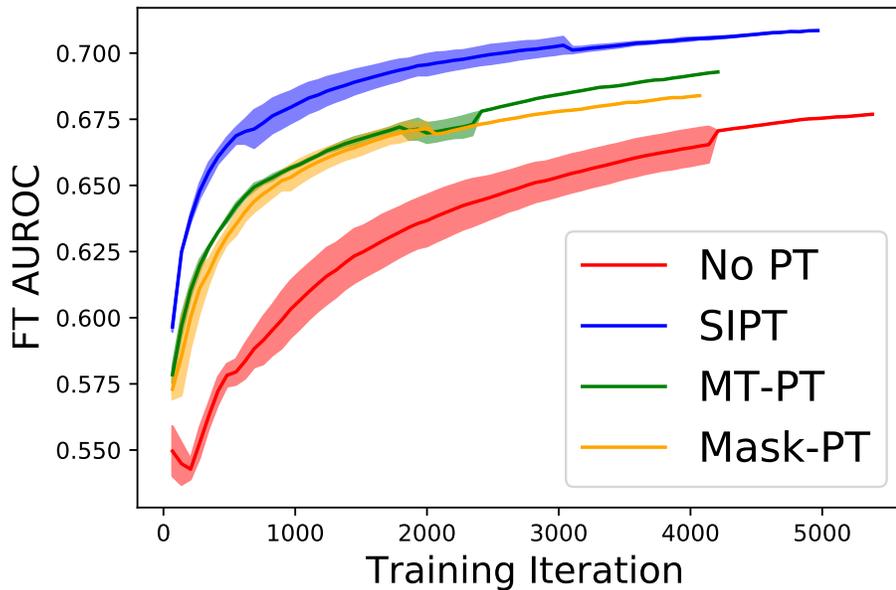


Figure 3: **Fine-tuning (FT) Performance over NETWORKS:** FT AUROC as a function of FT iteration for the NETWORKS dataset. The SIPT method converges faster and performs better than intra-sample (masked node modelling) or per-sample (multi-task classification) pre-training.

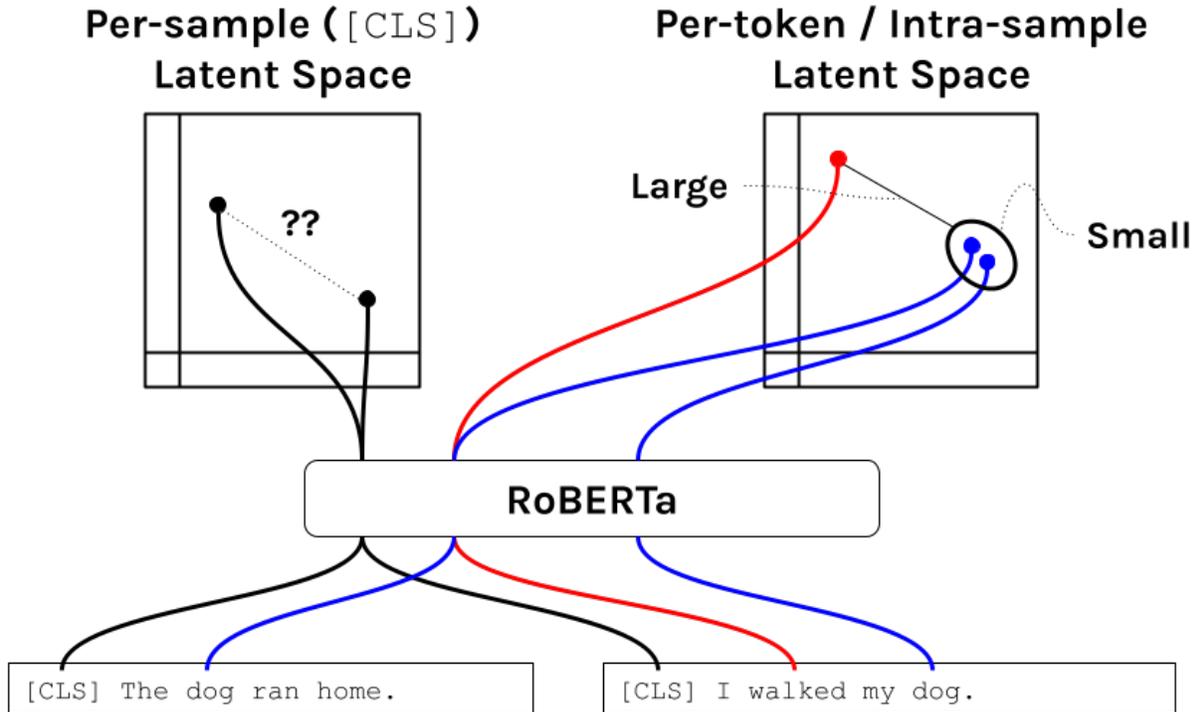


Figure 4: **Per-sample vs. Per-token Latent Space** Language model pre-training methods produce both per-sample and per-token latent spaces. Traditional language modelling objectives (illustrated here via the RoBERTa [4] model, which uses only a masked language model loss during pre-training) only constrain the per-token latent space.

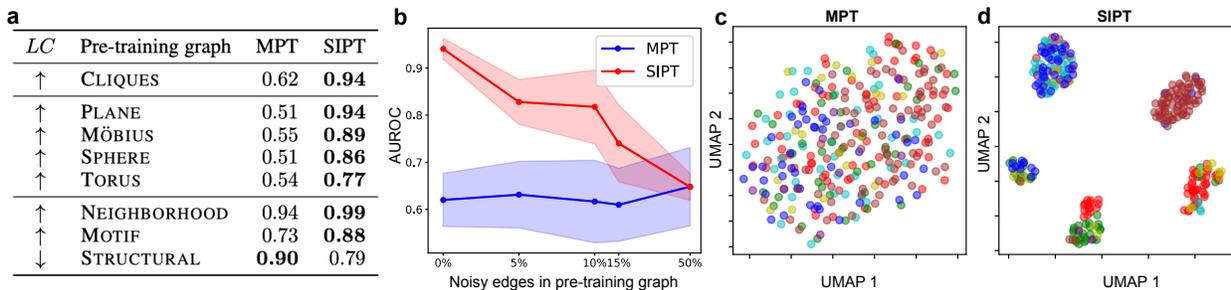


Figure 5: **Semi-synthetic Experiments Results:** (a) Comparisons between nearest-neighbor FT AUROC (higher is better) of LM PT models and SIPT models over various graphs with various forms of structural alignment. *LC* indicates the label consistency between FT task and  $G_{PT}$  (Definition 5). (b) Nearest-neighbor FT AUROC vs. noise rate. Up to 10% noise SIPT dramatically outperforms LM PT, and at 50% noise, the two approaches are equal. (c-d) Embedding space of MPT and SIPT models on the MÖBIUS dataset. Point colors indicate topic labels. SIPT’s embedding space reflects the structure of the PT graph, whereas MPT does not.

	PROTEINS	ABSTRACTS	NETWORKS
Data Modality ( $\mathbf{x}_i$ is a...)	Protein Sequence	Biomedical Paper Abstract	PPI Network Ego-graph
PT Dataset	Tree-of-life [74]	Microsoft Academic Graph [75, 76]	[43]
$G_{PT}$ : ( $\mathbf{x}_i, \mathbf{x}_j$ ) $\in E$ iff	$\mathbf{x}_i$ interacts with $\mathbf{x}_j$	$\mathbf{x}_i$ 's paper cites $\mathbf{x}_j$ 's paper	$\mathbf{x}_i$ 's central protein agrees on all but 9 Gene Ontology (GO) labels with $\mathbf{x}_j$ 's central protein.
Per-token baseline	TAPE [15]	SciBERT [91]	Attribute Masking [43]
Per-sample baseline	PLUS [45]	None	Multi-task learning [43]
FT Dataset	TAPE [15]	SciBERT [91]	[43]

**Table 1:** A summary of our datasets, tasks, and benchmarks. For example, for the PROTEINS domain, our pre-training dataset is the set of protein sequences contained in the tree-of-life dataset [74], proteins are linked in our pre-training graph  $G_{PT}$  if and only if they interact according to the tree-of-life graph, and we compare over the fine-tuning tasks in the TAPE benchmark against both the raw, per-token baseline publicly available in the TAPE model [15] as well as the per-sample baseline published in the PLUS pre-training model [45].

Domain	Task	Vs. Per-Token PT		vs. Per-Sample	
		RRE	$\Delta$	RRE	$\Delta$
PROTEINS	RH	<b>7.0%</b> $\pm$ 1.2	$\uparrow$	<b>8.4%</b> $\pm$ 2.4	$\uparrow$
	FL	-0.8%\pm1.3	$\sim$	<b>12.8%</b> $\pm$ 1.1	$\uparrow$
	ST	<b>13.1%</b> $\pm$ 2.5	$\uparrow$	2.2%\pm2.8	$\sim$
	SS	<b>4.5%</b> $\pm$ 0.2	$\uparrow$	<b>4.5%</b> $\pm$ 0.2	$\uparrow$
	CP	<b>10.5%</b> *	$\uparrow$	N/A	
ABSTRACTS	PF	0.3%\pm0.2	$\sim$	N/A	
	SC	2.4%\pm4.1	$\sim$	N/A	
	AA	<b>17.7%</b> $\pm$ 6.5	$\uparrow$	N/A	
	SRE	<b>6.7%</b> $\pm$ 0.4	$\uparrow$	N/A	
NETWORKS		7.8%\pm5.2	$\sim$	5.1%\pm2.7	$\uparrow$

**Table 2:** Relative reduction of error (RRE; defined to be  $\frac{[\text{baseline error}] - [G_{PT} \text{ model error}]}{[\text{baseline error}]}$ ) of models trained under our framework vs. published per-token or per-sample baselines. Higher numbers indicate models under our framework reduce error more and thus outperform baselines. The  $\Delta$  column indicates whether the model offers a statistically significant improvement ( $\uparrow$ ), no significant change ( $\sim$ ), or a statistically significant decrease ( $\downarrow$ ). Statistical significance is assessed via a  $t$ -test at significance level  $p < 0.1$ . Per-sample analysis and variance estimates for CP were infeasible due to the computational cost of this task.

FT Dataset	FT Task		Description	Metric
	Name	Abbr.		
TAPE [15]	Remote Homology	RH	Per-sequence classification task to predict protein fold category.	Accuracy
	Secondary Structure	SS	Per-token classification task to predict amino acid structural properties.	Accuracy
	Stability	ST	Per-sequence, regression task to predict stability.	Spearman’s $\rho$
	Fluorescence	FL	Per-sequence, regression task to predict fluorescence.	Spearman’s $\rho$
	Contact Prediction	CP	Intra-sequence classification to predict which pairs of amino acids are in contact in the protein’s 3D conformation.	Precision @ $L/5$
SciBERT [91]	Paper Field	PF	Per-sentence classification problem to predict a paper’s area of study from its title.	Macro-F1
	SciCite	SC	Per-sentence classification problem to predict citation intent	Macro-F1
	ACL-ARC	AA	Per-sentence classification problem to predict citation intent	Macro-F1
	SciERC	SRE	Per-sentence relation extraction	Macro-F1
NETWORKS [43]			Multi-label binary classification into 40 Gene Ontology terms	Macro-AUROC

**Table 3:** Fine-tuning tasks.

## Online Methods

### Per-token vs. Per-sample Latent Space: Definition of $\mathcal{Z}$

Let  $f_\theta$  be a pre-training (PT) model trained over a dataset  $\mathcal{X} \in \mathcal{X}^{N_{\text{PT}}}$ . Furthermore, let us assume that samples  $\mathbf{x} \in \mathcal{X}$  are composed of smaller units (*e.g.* tokens, sequence time-points, nodes in a network, etc.). Let us denote this by saying that  $\mathbf{x} = w_1, w_2, \dots, w_{n_x}$ . Finally, as is true in natural language processing (NLP) and NLP-derived settings, we assume that  $f_\theta$  can be seen to produce output embeddings for both the entire sample  $\mathbf{x}$ —which we will denote by  $f_\theta(\mathbf{x})$ —and for the internal tokens individually—which will denote by  $f_\theta(w_j|\mathbf{x})$ . For example, in the BERT model [35],  $f_\theta(\mathbf{x})$  will be given by the output embedding of the [CLS] token of  $\mathbf{x}$  and  $f_\theta(w_j|\mathbf{x})$  will be given by the output embedding of the  $j$ -th token in  $\mathbf{x}$ .

We can then formally define the per-sample latent space,  $\mathcal{Z}^{(S)}$  (which we will also refer to as  $\mathcal{Z}$  without the superscript), and the per-token (aka intra-sample) latent space  $\mathcal{Z}^{(T)}$  (Definitions 3 & 4, and Figure 4).

**Definition 3.** Per-Sample Latent Space We define the *per-sample latent space* induced by  $f_\theta$  as  $\mathcal{Z}^{(S)} = \{f_\theta(\mathbf{x})|\mathbf{x} \in \mathcal{X}\}$ . We will also use  $\mathcal{Z}$  with no superscript to refer to this space.

**Definition 4.** Per-token/Intra-sample Latent Space We define the *per-token latent space* (also known as the *intra-sample latent space*) induced by  $f_\theta$  as  $\mathcal{Z}^{(T)} = \{f_\theta(w_j|\mathbf{x})|w_j \in \mathbf{x}, \mathbf{x} \in \mathcal{X}\}$ .

Both of these spaces are very different and are useful in different contexts; for a task like named entity recognition, where the unit of classification is a single or short span of tokens, analyzing the per-token latent space will be more informative, whereas for a task like sentiment analysis, where the unit of classification is an entire sample (sentence), the per-sample latent space would be preferred [35]. Furthermore, another key difference between these spaces is that the traditional PT language model objective only induces significant constraints on the geometry of the per-token latent space and does not impact the per-sample latent space at all. This illustrates a gap in the capabilities of PT methods. In our work, we are concerned with precisely this gap and focus our attention on  $\mathcal{Z}$  (*i.e.*  $\mathcal{Z}^{(S)}$ ). We focus our attention on the per-sample latent space for 3 reasons:

1. There has been significantly more research on how to regularize the per-token latent space than the per-sample latent space, as we show in our extensive review (Table 4).
2. In many domains outside of NLP, the per-sample latent space is often of much greater interest than the intra-sample latent space. For example, in modelling protein sequences [15], drug structures [43], or electronic health record time series [46], per-sample tasks are of much greater interest than intra-sample tasks.
3. Even within NLP, modern methods struggle much more with representing whole passages

of text rather than short, isolated spans. This is evidenced by the battery of work examining sentence representations atop pre-trained language models [73, 92].

## Why is NLP Different than Other Domains?

In this work, we have implicitly argued that because a PT objective like masked language modelling (MLM) will not necessarily directly enrich the per-sample latent space  $\mathcal{Z}^{(S)}$ , it may yield models less well suited to downstream per-sample tasks than other approaches. One seeming contradiction to this is that methods in NLP like RoBERTa [4] (for which MLM is the only PT objective) succeed across both per-token and per-sample tasks.

In fact, this observation does not contradict our hypothesis but reflects a unique advantage of the natural language modality that does not apply in other domains. In particular, in the NLP domain (and not in other domains), we can leverage the flexibility of the language to sidestep any deficit in  $\mathcal{Z}^{(S)}$  by re-framing per-sample tasks as per-token, language modelling tasks. Significant literature exists documenting this phenomenon through the lenses of prompting, cloze-filling models, text-to-text transformers, and theoretical analyses [2, 3, 11, 77, 93]. For example, [93] examines the efficacy of pre-trained language models on sentiment analysis explicitly and show that the language modelling component alone can be used in a per-token manner to indirectly solve a review sentiment analysis task by judging the likelihood of following the review with a “:)” emoji vs. a “:(” emoji. In this way, they shift the *per-sample* task of sentiment analysis to a *per-token* task via the (inserted) emoji.

However, language model pre-training has also inspired many derived methods to be used in other non-NLP domains. For example, in modelling graphs, [43] has examined vertex or edge-masking strategies reminiscent of MLM, with vertices and edges analogous to tokens and entire graphs whole samples; in modelling time series data, [46] has examined masked imputation models, with timepoints analogous to tokens and whole time series to samples; and in modelling protein sequences, [45] has used masked language modelling directly, with individual amino acids representing tokens and entire proteins representing samples. *In all three of these domains, we cannot re-frame per-sample tasks as “per-token” tasks as we can in NLP, and accordingly, the problem of insufficient per-sample latent space regularization will likely be much more severe in these domains. Accordingly, existing work, including the three works referenced above, all find that augmenting the language model pre-training task with additional, per-sample level supervised tasks can be beneficial, or even absolutely essential, to improving performance [43, 45, 46, 94].*

## Pre-training Review Methodology

Papers were selected via a manual search of the natural language processing (NLP) and NLP-derived pre-training methods (*i.e.*, methods focused primarily on other domains or on multi-modal domains were excluded) via Google Scholar as well as by crawling through references of papers

already included. Citation counts for each work were obtained via Google Scholar on August 2nd, 2022. Publication date (used to calculate citations per month since publication date) was computed as the earlier of either (1) the paper’s venue-specific date of publication or (2) the first submission date to the arXiv or BioRxiv platform, as referenced via an exact title match. A manual review was done to classify how pre-training methods constrain latent space geometry and assign subjective, numerical “shallow-deep” and “explicit-implicit” axes scores. In total, over 90 methods were examined, of which 71 were suitable for inclusion in numerical review results (Figure 1 and Table 4). All methods considered are summarized and categorized (and reasons for exclusions are given) in Appendix A.

## Further Analysis of Reviewed Methods

This work has extensively examined how existing pre-training methods constrain the *per-sample* latent space. However, it is also worth examining how these methods constrain the per-token latent space to demonstrate the extent to which per-sample objectives are under-explored in current pre-training research. To that end, we break down all of the studies included in our review not only by how they constrain their per-sample latent spaces but also by how they constrain their per-token latent spaces (Table 4). These groupings are also done at a greater granularity than the previously examined categories to offer more insight into which methods use which techniques. We see that not only are there more types of per-token latent space constraints leveraged (10 vs. 7), but also methods consistently leverage a much greater diversity of per-token constraints vs. per-sample constraints (1.45 per-token constraints per method vs. 0.58 per-sample constraints, on average). We can further see from Figure 1 that the citation volume for works in this space is also heavily concentrated around methods that first employ no per-sample PT objective, followed by methods that only impose shallow, explicit methods, which further establishes this research gap.

Method	Per-token										Per-sample						
	Masked, discriminative, or standard language modelling	Template/prompt-style multi-task language model training	Concatenate related sentences together	Named entity masking	Relation masking	Per-token knowledge graph alignment	Named entity recognition and linking	(Unconstrained) attention over a KG	Joint token and entity embeddings	Syntactic Knowledge Distillation	Single-task classification	Multi-task classification	Whole-sample graph alignment	Per-sample augmentation-based contrastive alignment	Multi-lingual cross-sample contrastive alignment	Unsupervised clustering	Contextual autoencoding
[1] ELMO	✓																
[2] GPT-3	✓																
[3] T5	✓	✓															
[4] RoBERTa	✓																
[5] GPT-1	✓																
[6] GPT-2	✓																
[7] BART	✓																
[8] Unsupervised Cross Lingual	✓																
[9] ELECTRA	✓																
[10] SpanBERT	✓																
[11] UniLM	✓																
[12] DAPT	✓																
[13] ERNIE (Sun et. al.)	✓			✓													
[14] KnowBERT	✓						✓	✓	✓								
[15] TAPE	✓																
[16] LUKE	✓			✓													
[17] Topp	✓	✓															
[18] Pretrained Encyclopedia	✓			✓													
[19] MSA	✓		✓														
[20] COLAKE	✓			✓	✓												
[21] BERTMK	✓																
[22] ERICA	✓					✓											
[23] JAKET	✓																
[24] CALM	✓	✓															
[25] KeBioLM	✓						✓	✓	✓								
[26] MG-BERT (Molecules)	✓																
[27] CDLM	✓		✓														
[28] KgPLM	✓			✓													
[29] kNN PT	✓		✓														
[30] LP-BERT	✓			✓	✓												
[31] MG-BERT (NLP)	✓																
[32] UD-PrLM	✓					✓											
[33] ESM-1B	✓																
[34] UniRep	✓																
[35] BERT	✓																
[36] ERNIE (Zhang et. al.)	✓						✓										
[37] CokeBERT	✓																
[38] SPIDER	✓					✓											
[39] Syntactic-Distilled BERT	✓								✓								
[40] ALBERT	✓																
[41] SMedBERT	✓					✓											
[42] MT-DNN	✓																
[43] Graph-PT	✓																
[44] SentiLARE	✓																
[45] PLUS	✓							✓									
[46] EHR-PT	✓																
[47] ERNIE 2.0 (Sun et. al.)	✓					✓											
[48] ERNIE 3.0 (Sun et. al.)	✓			✓	✓												
[49] Dict-BERT	✓		✓				✓										
[50] LinkBERT	✓																
[51] StructBERT	✓																
[52] MARGE	✓																
[53] REALM	✓	✓	✓	✓											✓		
[54] GraphCL	✓																
[55] GCC	✓																
[56] DeCLUTR	✓																
[57] CLEAR	✓																
[58] JOAO	✓																
[59] COCO-LM	✓																
[60] InfoWord	✓																
[61] MICRO-Graph	✓														✓		
[62] STS-CT	✓																
[63] CAPT	✓																
[64] GearNet	✓					✓											
[65] InfoXLM	✓													✓			
[66] GLM	✓																
[67] KCL	✓			✓													
[68] KEPLER	✓																
[69] CK-GNN	✓																
[70] XLM-K	✓																
[71] Webformer	✓																

**Table 4: Existing Pre-training (PT) Methods:** A subset of existing PT methods, broken down by how they constrain per-token and per-sample latent space geometries.

## Constraints on $\mathcal{L}_{\text{SI}}$ in our Framework

Formally, for  $\mathcal{L}_{\text{SI}}$  to be valid, then there must exist a distance function  $d : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ , radius  $r \in \mathbb{R}$ , and loss value  $\ell^* \in \mathbb{R}$  such that at any solution  $\theta^*$  for which  $\mathcal{L}_{\text{SI}}(\theta^*) < \ell^*$ , the learned embeddings  $z_i = f_{\theta^*}(\mathbf{x}_i)$  must recover the graph  $G_{\text{PT}}$  under a radius nearest neighbor connectivity algorithm via distance function  $d$  and radius  $r$ —*i.e.*, it must be the case that  $(\mathbf{x}_i, \mathbf{x}_j) \in E$  if and only if  $d(f_{\theta^*}(\mathbf{x}_i), f_{\theta^*}(\mathbf{x}_j)) < r$ . Furthermore, for the particular graph  $G_{\text{PT}}$  and latent space  $\mathcal{Z}$ , the set of  $\theta^*$  such that  $\mathcal{L}_{\text{SI}}(\theta^*) < \ell^*$  must be non-empty (*i.e.* such a solution must exist).

## Realizing Existing Methods in our Framework

Let  $\mathbf{X} \in \mathcal{X}^{\text{NPT}}$  be the pre-training dataset throughout this section. In cases where we have some auxiliary information (*e.g.*, supervised, per-sample, pre-training labels), they will be denoted by  $\mathbf{Y} \in \mathcal{Y}^{\text{NPT}}$ .

### Methods with no per-sample objectives

Naturally, we can realize any method that only employs a per-token pre-training objective within our framework simply by setting  $\lambda_{\text{SI}} = 0$ . This realization is trivial and offers no insight into the suitability of these pre-training methods for downstream per-sample tasks.

### Methods with a supervised, single-task per-sample objective (*e.g.*, BERT [35])

A simple, single-task, per-sample, classification pre-training objective induces a geometric constraint in the output latent space on the basis of the inner product “distance” between samples of the same vs. different class labels. We can use this observation to realize a reduction from a valid SIPT objective to the original classification objective. In particular, we can introduce a graph  $G = (\{\mathbf{x}_i \in \mathbf{X}\}, \{(\mathbf{x}_i, \mathbf{x}_j) | y_i = y_j\})$  which consists of cliques corresponding to each unique label  $c \in \mathcal{Y}$ . Then, leveraging any structure-preserving metric learning loss with a cosine distance objective will, at optimality, recover a solution that also satisfies the original classification objective, where we use centroids of the induced clique embeddings to represent class embeddings.

### Methods with a supervised, multi-task per-sample objective (*e.g.*, MT-DNN [42])

A slightly more complicated case is when methods employ a multi-task, per-sample classification objective. In this case, there are two ways to realize this task within the SIPT framework. First, we can simply transform the multi-task objective into a single-task objective by constructing a new label-space consisting of the Cartesian product of all label spaces for each task individually. This will greatly increase the number of “labels” in the task, but then the problem can be realized via a graph of disconnected cliques much like in the single-task setting.

However, there is another manner in which we can realize this objective in the SIPT framework; In particular, suppose our collection of tasks consists of  $k$  label spaces:  $\mathcal{Y} = \mathcal{Y}_1 \times \cdots \times \mathcal{Y}_k$ . Then, we can construct a graph  $G = (V, E)$  such that:

1. the vertices consist of all pre-training samples  $\mathbf{x}_i$  as well as auxiliary nodes corresponding to each label  $c_h^{(j)} \in \mathcal{Y}_h$  across each task:  $V = \{\mathbf{x}_i \in \mathbf{X}\} \cup \mathcal{Y}_1 \cup \dots \cup \mathcal{Y}_k$
2. the edges contain links between each sample  $\mathbf{x}_i$  and label  $y_h^{(i)}$  across all tasks  $1 \leq h \leq k$ :  $E = \{(\mathbf{x}_i, c_h^{(j)}) | y_h^{(i)} = c_h^{(j)}\}$ .

Then, we can see that if we solve the SIPT problem under a structure-preserving metric learning loss, we will naturally have produced embeddings for each  $\mathbf{x}_i$  which are close (in inner-product distance space) to the class embeddings corresponding to their labels for each task, while they are also far from other, non-matching class embeddings, as desired. This second approach is more useful to us in considering the ramifications of this style of constraint because it enables us to make more rigid theoretical guarantees via the SIPT theory.

### Methods with a based contrastive per-sample objective (e.g., GraphCL [54])

It is challenging to realize contrastive learning approaches within the SIPT framework, but it is still possible. Here, we highlight two distinct types of contrastive learning approaches we can capture within SIPT: a noising/augmentation-based approach, in which sample embeddings are constrained to be similar to embeddings of noised versions of said samples; and a multi-modal (or multi-lingual) contrastive approach, in which there exists a 1:1 mapping between two different sub-modalities within  $\mathbf{X}$  which is used to join those two modalities into a unified latent space (e.g. a model which constrains embeddings of English sentences to be close to embeddings of their french translations, but far from unrelated sentences).

To consider the augmentation/noising policy type first, let  $h : \mathbf{x}_i \mapsto \tilde{\mathbf{x}}_i$  represent the noising transformation. Then, to build an analogous SIPT model to this model, we construct an augmented dataset consisting of all original data points alongside all possible transformed versions of the original data points under  $h$ :  $\mathbf{X}' = \mathbf{X} \cup \left( \bigcup_{i=1}^{N_{PT}} \text{Im}(h|_{\mathbf{x}_i}) \right)$ . Note that even in contexts where  $h$  is continuous (and thus has an infinite image), we can still construct this dataset in practice because training is only performed over a finite number of steps, meaning our augmented dataset  $\mathbf{X}'$  need only be expanded to cover a finite number of augmentations. Then, the associated pre-training graph is simple; we simply use every sample in the augmented dataset  $\mathbf{X}'$  as a vertex and connect any two samples if and only if one is a transformed version of the other. This forms a graph of many disconnected stars (one star for each original datapoint  $\mathbf{x}_i$ ), and thus it does not directly enforce any particular geometry via our current theory. However, in cases where dataset size is sufficiently large,  $h$  sufficiently expressive, and data density sufficiently high, then the natural continuity of any neural network model will induce additional, auxiliary connections across these stars (if, for example, the noised versions of two distinct samples have a high probability of being very similar), which increases the depth of the geometric constraints enforced. Quantifying the exact parameters of these interactions, however, we leave to future work.

In the case of the multi-modal/multi-lingual contrastive alignment objective across  $k$  modalities, our setup is much simpler: we simply let  $G_{PT}$  be a  $k$ -partite graph whose samples consist of individual data points (across all modalities) and edges connect samples that compose a matching pair across modalities (*e.g.* edges link English sentences to their french translations). The extent to which this constrains the output geometry in practice, then, comes down to several questions: (1) Is the cross-modal alignment a one-to-one, one-to-many, or many-to-many alignment (which impacts the geometry of the resulting graph), (2) How large and dense is the dataset (which impacts the extent to which additional, indirect edges will be induced due to continuity in practice), and (3) How do other pre-training objectives constrain the individual modalities separately? In a case where this graph is one-to-one, and no other constraints are induced in each modality separately, this objective will offer only minimal constraints as the resulting graph will consist of many disconnected 2-cliques.

### Methods with a per-sample graph-alignment objective (*e.g.*, KEPLER [68])

Methods that explicitly align samples with a provided pre-training graph (KEPLER [68], CK-GNN [69], XLM-K [70], and WebFormer [71]) are naturally already realized within SIPT, so need no further commentary here.

## Structure-inducing Losses Examined in this Study

### Multi-similarity loss

The multi-similarity loss, parametrized by  $w_+$ ,  $w_-$ , and  $t$ , is given below:

$$\mathcal{L}_{SI} = \frac{1}{Nw_+} \log \left( 1 + \sum_{(i,j) \in E} e^{-w_+ (\langle f_{\theta}(\mathbf{x}_i), f_{\theta}(\mathbf{x}_j) \rangle - t)} \right) + \frac{1}{Nw_-} \log \left( 1 + \sum_{(i,j) \notin E} e^{w_- (\langle f_{\theta}(\mathbf{x}_i), f_{\theta}(\mathbf{x}_j) \rangle - t)} \right),$$

### Contrastive loss

Our contrastive loss is modeled after [89]’s version. For this loss, we assume we are given the following mappings: ‘pos’, which maps  $\mathbf{x}$  into a positive node (*i.e.*, linked to  $\mathbf{x}$  in  $G_{PT}$ ), and ‘neg’, which maps  $\mathbf{x}$  into a negative node (*i.e.*, not linked to  $\mathbf{x}$  in  $G_{PT}$ ). The union of a seed minibatch  $B$  of points  $\mathbf{X}_B$  and its images under ‘pos’ and ‘neg’ mappings form a full minibatch. This loss is specified by the positive and negative margin parameters  $\mu_+$  and  $\mu_-$  as:

$$\mathcal{L}_{SI}^{(CL)} = \frac{1}{N} \sum_{\mathbf{x}_i \in \mathbf{X}} \max(\mathcal{D}(\mathbf{x}_i, \text{pos}(\mathbf{x}_i)) - \mu_+, 0) + \frac{1}{N} \sum_{\mathbf{x}_i \in \mathbf{X}} \max(\mu_- - \mathcal{D}(\mathbf{x}_i, \text{neg}(\mathbf{x}_i)), 0).$$

### Additional Choices within the SIPT Framework

In addition to a loss term, we can use negative sampling to improve efficiency. Using the full graph  $G_{PT}$ , which is not available in many contexts where negative sampling is employed, we

can leverage the distance between samples calculated on  $G_{\text{PT}}$ , which provides a complementary source of information beyond embedding space distance alone. For example, one could use this to limit negative samples within the same connected component, but more complex strategies based on graph sampling (e.g. [95]) could also be used.

## Proof of Theorem 1

We begin by defining the notion of “Local Consistency,” which (informally) quantifies how “smooth” a given fine-tuning task label is over a graph  $G_{\text{PT}}$  (Definition 5). In addition, note that throughout all proofs, we will assume that the PT and FT datasets are iid, that FT tasks, though they may be unobserved over PT samples, are well defined over the entire PT and FT domain and thus true labels do exist (though they may be unknown) for PT samples, and that the sampling distribution of the PT/FT data has full support over the label-space of any considered task.

**Definition 5** (Local Consistency). Let  $y : X \rightarrow \mathcal{Y}$  be a task over a domain  $X$ , and let  $G = (V, E)$  be a graph such that  $X \subseteq V$ . The *local consistency*  $\text{LC}_G(y)$  is the probability that a node’s label  $y(x)$  agrees with the majority of labels of  $x$ ’s neighbors in  $G$ :

$$\text{LC}_G(y) = \mathbb{P} \left( y(x) = \underset{c \in \mathcal{Y}}{\operatorname{argmax}} \sum_{x' \in X | (x, x') \in E} \mathbb{1}_{y(x')=c} \right).$$

Note this is closely related to *homophily* [96–98].

With Local Consistency defined, we can now formally prove Theorem 1, reproduced below.

**Theorem 1.** Let  $\mathbf{X}_{\text{PT}}$  be a PT dataset,  $G_{\text{PT}}$  be a PT graph, and let  $f_{\theta^*}$  be an encoder pre-trained under a PT objective permissible under our framing that realizes a  $\mathcal{L}_{\text{SI}}$  value no more than  $\ell^*$ . Then, under embedder  $f$ , the nearest-neighbor accuracy for a FT task  $y$  converges as dataset size increases to at least the local consistency (Definition 5) of  $y$  over  $G_{\text{PT}}$ .

*Proof.* Given  $f$  realizes SIPT-optimal embeddings, we know that if we define a  $r$ -NN predictor via the same radius  $r^*$  at which  $f$  achieves optimality, then this predictor will be correct exactly as often as the label of a given node in the graph  $G_{\text{PT}}$  agrees with the labels of its neighbors—which is  $\text{LC}_{G_{\text{PT}}}(y)$ . This classifier may not be well defined for small FT dataset sizes. However, as if data is not sufficiently dense, there may be no data points within the radius  $r$  of a given query. Similarly, without sufficient PT data, the LC computed over the empirical distribution of the graph  $G_{\text{PT}}$  may be a poor proxy for the true distribution. As PT and FT dataset sizes increase, however, we can achieve at least this performance. We may be able to achieve even higher performance if other effects motivate stronger performance at radii smaller than  $r^*$ , but this is not guaranteed.  $\square$

## Proof of Corollary 1

**Corollary 1.** Let  $\mathbf{X}_{\text{PT}} \in \mathcal{X}^N$ , be a PT dataset with corresponding labels  $\mathbf{y} \in \mathcal{Y}_{\text{PT}}^N$ . Define  $G_{\text{PT}} = (\mathbf{X}_{\text{PT}}, E)$  such that  $(\mathbf{x}_i, \mathbf{x}_j) \in E$  if and only if  $y_i = y_j$ .

Then, the local consistency for a given FT task  $\mathbf{y}^{(\text{FT})}$  over  $G_{\text{PT}}$  (and thus by Theorem 1, the nearest-neighbor accuracy for any optimized SIPT embedder) is upper bounded by the probability that a sample  $x_i$ 's fine-tuning label  $y_i^{(\text{FT})}$  agrees with the majority class label for task  $\mathbf{y}^{(\text{FT})}$  over the clique consisting of all nodes with the same pre-training label  $y_i$  as  $x_i$ .

*Proof.* This follows directly from the definition of Local Consistency,  $G_{\text{PT}}$ , and the law of total probability. In particular,

$$\begin{aligned} \text{LC}_{G_{\text{PT}}}(y_{\text{FT}}) &= \mathbb{P} \left( y_{\text{FT}}(\mathbf{x}_i) = \underset{\ell \in \mathcal{Y}_{\text{FT}}}{\text{argmax}} \sum_{\mathbf{x}_j \in \mathbf{X}_{\text{PT}} | (\mathbf{x}_i, \mathbf{x}_j) \in E(G_{\text{PT}})} \mathbb{1}_{y_{\text{FT}}(\mathbf{x}_i) = \ell} \right) \\ &= \mathbb{P}(y_{\text{FT}}(\mathbf{x}_i) = \text{MC}(\mathbf{x}_i, y_{\text{FT}})) \\ &= \sum_{\ell_{\text{PT}} \in \mathcal{Y}_{\text{PT}}} \mathbb{P}(y_i = \ell_{\text{PT}}) \mathbb{P}(y_{\text{FT}}(\mathbf{x}_i) = \text{MC}(\mathbf{x}_i, y_{\text{FT}}) | y_i = \ell), \end{aligned}$$

With Local consistency found, a simple application of Theorem 1 completes the proof.  $\square$

Note that this has a dependence on the PT dataset size as the probabilities  $\mathbb{P}$  are taken over the empirical distribution induced by the dataset  $\mathbf{X}_{\text{PT}}$  and graph  $G_{\text{PT}}$  inherent in local consistency — if  $\mathbf{X}_{\text{PT}}$  is too small, these empirical distributions will be poor proxies for the true distribution and this bound will not hold tightly. However, once saturation is reached, it will not improve beyond this fixed upper bound relating to task correlation.

## Proof of Corollary 2

**Corollary 2.** Let  $\mathbf{X}_{\text{PT}}$  be a PT dataset that can be realized over a valid manifold  $\mathcal{M}$ . Assume  $\mathbf{X}_{\text{PT}}$  is sampled with full support over  $\mathcal{M}$ . Let  $G_{\text{PT}}(\mathbf{X}_{\text{PT}}, E)$  be an  $r$ -nearest-neighbor graph over  $\mathcal{M}$  (e.g.,  $(\mathbf{x}_i, \mathbf{x}_j) \in E$  if and only if the geodesic distance between the two points on  $\mathcal{M}$  is less than  $r$ :  $\mathcal{D}_{\mathcal{M}}(\mathbf{x}_i, \mathbf{x}_j) < r$ ). Let  $\mathbf{y}^{(\text{FT})}$  be a FT classification task that is almost everywhere smooth on the manifold.

Then, as PT dataset size (and thus the size of  $G_{\text{PT}}$ ) tends to  $\infty$ , and  $r$  tends to zero, the local consistency of  $\mathbf{y}^{(\text{FT})}$  over  $G_{\text{PT}}$  (and thus by Theorem 1 the nearest-neighbor accuracy of an SIPT embedder) will likewise tend to 1.

*Proof.* As  $r \rightarrow 0$ , provided PT dataset size increases at a sufficient associated rate so as to maintain a constant minimum degree of  $G$ , we have the property that the total diameter over  $\mathcal{M}$  contained in a node's local neighborhood within  $G_{\text{PT}}$  likewise decreases. Given some fixed node  $\mathbf{x} \in \mathcal{M}$

that is within the interior of a set of constant  $y_{\text{FT}}$  label, this implies that, eventually, it will grow sufficiently small that all of  $\mathbf{x}$ 's neighbors share the same label as  $\mathbf{x}$  under  $y_{\text{FT}}$ .

More concretely, it is clear that this point will occur exactly when  $r$  is the geodesic distance between  $\mathbf{x}$  and the boundary of the surrounding constant-label patch containing  $\mathbf{x}$ . But, it is clear that the only sections of  $\mathcal{M}$  will not have the property that neighborhoods around points will be constant w.r.t.  $y_{\text{FT}}$  labels will almost everywhere be patches within distance  $r$  of the points where  $y_{\text{FT}}$  changes.

This implies that as  $r \rightarrow 0$ , then almost everywhere will the neighborhoods around a node  $\mathbf{x}$  be constant w.r.t.  $y_{\text{FT}}$ . However, this implies that almost everywhere would  $y_{\text{FT}}$  display perfect local consistency, as desired.  $\square$

## Semi-synthetic Experiments Validating Theoretical Results

We can further validate the theoretical analyses of our framework via semi-synthetic experiments. In particular, we create several datasets of natural language sentences augmented with synthetic graphs with known relationships to certain FT tasks (e.g., low or high local consistency, low or high rates of noise). We then use these datasets to validate three important properties of PT methods: First, do PT methods trained with a  $\mathcal{L}_{\text{SI}}$  and  $G_{\text{PT}}$  yield Nearest-neighbor FT performance in accordance with our theory? In particular, do (a) FT tasks with high local consistency over the PT graph offer better performance, and (b) those with very low local consistency offer worse performance? Second, do PT methods trained with a  $\mathcal{L}_{\text{SI}}$  and  $G_{\text{PT}}$  suffer significantly when pre-training graphs are polluted with noise? Finally, third, do the latent space geometry regularizing properties of  $\mathcal{L}_{\text{SI}}$  yield methods whose embeddings more clearly cluster than embeddings produced by traditional pre-training alone?

### Pre-training & fine-tuning datasets

Across all experiments, our synthetic datasets consist of free-text sentences from <https://www.kaggle.com/mikeortman/wikipedia-sentences> (CC BY-SA 4.0 License).

Topics were assigned to these sentences by running Latent Dirichlet Allocation via Scikit-learn [99] over a Bag-of-words representation to 100 topics, with otherwise default parameters. Given the probabilities over all 100 topics, we treated the prediction of the most probable topic as a 100-class multi-class classification problem for our FT task in these experiments.

To test across various graphs, we produce a number of pre-training graphs per experiment, as detailed below.

### Pre-training graphs

We use graphs spanning 3 categories. (1) A graph (CLIQUEs) consisting of disconnected cliques, where sentences are linked in the graph if they share the same topic label. (2) Graphs composed of nearest-neighbor graphs defined over simplicial manifolds built using topic probabilities to lo-

calize sentences onto simplices. We explore manifolds with a range of topological complexity, including: PLANE, MÖBIUS, SPHERE, and TORUS. Finally, (3) we define three graphs according to a mechanistic process that allows us to control how topic labels relate to graph structure: first, so that topics are maximally conserved within local neighborhoods (NEIGHBORHOOD); second, by assigning sentences to nodes in the graph such that each graph motif corresponds to a unique topic (MOTIF); and third, such that node topics are driven by non-local graph structural features, on the basis of graphlet degree vectors (STRUCTURAL). Details for each pre-training graph formation are given below.

### CLIQUEs Graph Setup

To construct the Cliques graph setting, we choose a random subset of sentences as  $\mathbf{X}_{PT}$  and define  $G_{PT} = (\mathbf{X}_{PT}, E)$  such that  $(x_i, x_j) \in E$  if and only if  $x_i$  and  $x_j$  share the same topic label.

### PLANE, MÖBIUS, SPHERE, & TORUS Graphs

For these graphs, we take a more involved practice to localize sentences onto specifiable simplicial manifolds, then construct pre-training graphs via radius nearest neighbor graphs on those manifolds. This involves several steps:

**Localizing Sentences on Simplices** We can localize any sentence in our overall dataset onto a 2-simplex by mapping them onto the (re-normalized) probabilities associated with their top-3 topics. Doing this means that the simplex on which they are localized has vertices corresponding to possible topics among our 100 total topics.

**Stitching Topic-simplices Into Manifolds** Given these topic-simplex localized sentences, we need to construct our manifolds. To do so, we first produce any arbitrary simplicial tiling of a 2-manifold. With this tiling, all that remains to localize sentences onto the manifold is to find a self-consistent mapping of topics to simplex vertices (in the tiling) such that all topic-simplices induced by this mapping have sufficiently many associated samples to enable roughly uniform sampling.

**Sampling Points** After finding a self-consistent map of topics to simplicial tiling vertices that satisfy density requirements, we can sample sentences onto the manifold. To make this process more uniform, we also calculate the relative entropy of each sentence (over the re-normalized probabilities of the top-3 topics), bin those entropies into buckets, then sample first what entropy bucket we wish to draw from such that the induced distribution of sentence entropies is approximately uniform, then sample within that entropy bucket.

**Calculating on-Manifold Distances** Finally, with sentences sampled and localized onto a simplicial manifold, we then need to compute approximate geodesic distances to enable building

radius-nearest-neighbor graphs over these sentences. To do so, we use an approximate algorithm that considers only on-simplex distance (*e.g.*, it does not consider any curvature penalties) which is equivalent to calculating the distance between any pair of points over the simplices presuming they were flattened onto a plane (this flattening naturally does not preserve manifold topology, but along only the shortest path between any particular set of two points it is always possible to do so with a 2-manifold).

The above process describes how to produce a radius-nearest-neighbor graph for any specifiable manifold using our topic-model outputs. We do this for simplicial manifolds that correspond topologically to a simple plane (PLANE), a möbius strip (MÖBIUS), a sphere (SPHERE), and a torus (TORUS).

### **STRUCTURAL, NEIGHBORHOOD & MOTIFS Graphs**

In order to form these examples, we must (1) define our overall graphs, (2) featurize these graphs in a manner that is reflective of different forms of graph structure, then (3) use these featurizations to assign sentences to graph nodes to form our pre-training dataset.

**Graph Construction** We sample graphs by first building a base cycle of a parametrized size, then add motifs along this cycle by sampling small graphs from all possible connected graphs of size less than 6 nodes.

**Node Featurization** Nodes in this graph are then assigned internal features based on three notions of graph topology. For the “Neighborhood” label, a node  $n$  is identified according to an index-vector indicating which nodes in the graph are within shortest-path distance 3 of  $n$ . For the “Motif” label,  $n$  is identified based on its membership either in the base cycle or any of the attached random subgraphs. For the “Structural” label,  $n$  is identified based on its graphlet degree vector (of order 4). For structural and homophily features, categorical labels are then produced by feeding these raw representations through a  $k$ -means clustering algorithm.

**Sentence Assignment** We assign sentences to nodes in multiple ways so that we can produce datasets that reflect each of the notions of graph structure discussed previously. In particular, for either the neighborhood, motif, or structural labels, each sentence topic is matched to a node label, then sentences are assigned randomly to nodes in the graph with a matching topic label. Note that this produces a dataset where the graph structure is only partially reflected by the node’s features, which is itself another useful test of the SIPT method, as it would not be useful if SIPT could only capture data in contexts where the graph was perfectly reflected by the node features themselves.

### **Expected local consistency between graphs $G_{PT}$ and the topic prediction FT task**

Of all these graphs, we expect that topics will display a low local consistency over the STRUCTURAL graph and a moderately high local consistency over the MOTIF graph (as graph motifs are all connected components), and high local consistency everywhere else.

### **Network Architecture & Hyperparameters**

The Cliques and Mechanistic experiments use a shallow Transformer model with 2 layers and 10 hidden units. The Manifold experiments use a 3-layer Transformer model with 256 hidden units. Hyperparameters were not tuned but were chosen by hand to produce as small a network as possible while permitting reasonable learning dynamics.

### **Experimental setup**

To answer our three questions, we will pre-train models under both traditional LM pre-training alone and a new, structure-inducing PT (SIPT) method within our paradigm that augments the loss with a contrastive learning loss over  $G_{PT}$ , with  $\lambda_{SI} = 0.1$ . Both models use a shallow transformer encoder for  $f_\theta$  and a character-level tokenization scheme. Final results are reported via the AUROC of 3-nearest-neighbor classifiers over the latent space, per-sample embeddings. In line with our theoretical predictions, we expect to see higher NN FT performance in all settings where the FT task (topic prediction) has high local consistency over the graph  $G_{PT}$  (all graphs except STRUCTURAL) and worse performance in the case where the local consistency is very low (STRUCTURAL).

We also assess the stability of our method as the graph  $G_{PT}$  is noised using the CLIQUES graph by randomly adding additional edges with varying rates.

### **Semi-synthetic Result 1: SIPT improves performance over LM PT by $0.26 \pm 0.13$ AUROC on graphs where the topic task has a high local consistency**

As can be seen in Figure 5a, SIPT offers significant improvements over LM PT in nearest-neighbor FT AUROC across all graph types with strong topic local consistency.

### **Semi-synthetic Result 2: SIPT’s empirical results are in agreement with theoretical findings**

In line with our theoretical findings, SIPT only under-performs LM PT on the STRUCTURAL graph where the topic task (by design) does not have strong local consistency. This validates our theoretical results by showing that local consistency strongly predicts Nearest-neighbor FT performance.

### **Semi-synthetic Result 3: SIPT is robust to incomplete and noisy pre-training graphs**

Figure 5b shows Nearest-neighbor FT AUROC as a function of noise rate on the CLIQUES graph. For up to 15% noise, SIPT shows improvements over LM PT, and even at 50% noise, the two approaches perform comparably.

## Semi-synthetic Result 4: SIPT pre-trained embeddings show stronger clustering than LM PT embeddings

Figure 5c-d shows embeddings produced under the MÖBIUS graph either by LM PT or SIPT, clustered via UMAP into 2 dimensions. It is clear visually from these figures that SIPT embeddings show clear clusters strongly associated with the topic-modelling FT task, whereas LM PT embeddings do not.

### Conclusions

From these analyses, we see that augmenting PT with per-sample structure-inducing objectives can both (1) offer significant advantages over existing PT architectures and (2) permit analytical reasoning about which FT tasks PT will offer improvements. These findings are not surprising; in these semi-synthetic experiments, we designed our graphs explicitly to have either high or low local consistency with respect to our FT task so that we could probe exactly whether SIPT methods would behave in accordance with theory in tightly controlled settings. In this way, the graphs  $G_{PT}$  used here may not be reflective of graphs in the real world, which will be chosen more independently of specific FT tasks. To address this, in the Results section, we demonstrate experimental results over diverse real-world datasets with real, FT-task-independent graphs to show that the gains persist in more realistic scenarios.

## Further Details on Real-world Experiments

### Further Details on the PROTEINS Dataset and FT tasks

**PT Dataset** We use a dataset of  $\sim 1.5M$  protein sequences from the Stanford Tree-of-life dataset [74] (<https://snap.stanford.edu/tree-of-life/data.html>). The associated Github repository for this resource lists an MIT license.

**PT Graph** Two proteins are linked in  $G_{PT}$  if and only if they are documented in the scientific literature to interact, according to the tree-of-life interaction dataset. This is an external knowledge graph.

**FT Dataset/Tasks** We use the TAPE FT benchmark tasks [15], including Remote homology (RH), a per-sequence classification task to predict protein fold category (metric: accuracy); Secondary structure (SS), a per-token classification task to predict amino acid structural properties (metric: accuracy); Stability (ST) & Fluorescence (FL), per-sequence, regression tasks to predict a protein’s stability and fluorescence, respectively (metric: Spearman’s  $\rho$ ); and Contact prediction (CP), an intra-sequence classification task to predict which pairs of amino acids are in contact in the protein’s 3D conformation (metric: Precision at  $L/5$ ).

**Baselines** We compare against the published TAPE model [15], which uses an LM task alone as our per-token comparison point, and the PLUS [45] model, which optimizes for LM and

supervised classification jointly, for our per-sample comparison point.

The tasks in the TAPE benchmark [15] on which we test are described more fully below. All these datasets are publicly available. All datasets can be obtained directly on TAPE’s Github (<https://github.com/songlab-cal/tape#data>), which lists no licenses for these datasets though the overall Github is released under a BSD 3-Clause ”New” or ”Revised” License.

**Remote Homology** This is a per-sequence, multi-class classification problem, evaluated using accuracy, which tasks a model to predict a protein fold category at a per-sequence level. This task’s dataset contains 12,312/736/718 train/val/test proteins and is originally sourced from [100].

**Secondary Structure** This is a per-token, multi-class classification problem, evaluated using accuracy, which tasks a model to predict the structural properties of each amino acid in the final, folded protein. This task’s dataset contains 8,678/2,170/513 train/val/test proteins, and is originally sourced from [101].

**Stability** This is a per-sequence, continuous regression problem evaluated using the Spearman correlation coefficient, which tasks a model to predict the protein’s stability in response to environmental conditions. This task’s dataset contains 53,679/2,447/12,839 train/val/test proteins, and is originally sourced from [102].

**Fluorescence** This is a per-sequence, continuous regression problem evaluated using the Spearman correlation coefficient, which tasks a model to predict how brightly a protein will fluoresce. This task’s dataset contains 21,446/5,362/27,217 train/val/test proteins, and is originally sourced from [103].

### **Further Details on the ABSTRACTS Dataset and FT tasks**

**PT Dataset** We use a dataset of  $\sim 650\text{K}$  free-text scientific article abstracts from the Microsoft Academic Graph (MAG) dataset [75, 76]. The ABSTRACTS PT data (the Microsoft Academic Graph dataset) is licensed with an Open Data Commons Attribution License (ODC-By) v1.0 license.

**PT Graph** Two abstracts are linked in  $G_{PT}$  if and only if their corresponding papers cite one another. This is a self-supervised graph.

**FT Dataset/Task** We use a subset of the fine-tuning tasks used in the SciBERT paper [91], including Paper field (PF), SciCite (SC), ACL-ARC (AA), and SciERC Relation Extraction (SRE), all of which are per-sentence classification problems (metric: Macro-F1). PF tasks models to predict a paper’s area of study from its title, SC & AA tasks both predict an “intent” label for citations, and SRE is a relation extraction task.

**Baseline** We compare against the published SciBERT model [91] as our per-token comparison and lack an associated per-sample comparison as we don’t know of any published per-sample models in the academic papers modality.

The tasks in the SciBERT benchmark [91] on which we test are described more fully below. All tasks here are per-sentence, multi-class classification problems (i.e., we do not study any per-token tasks), and all are evaluated in Macro-F1 (out of 1). All FT datasets can be obtained from the SciBERT Github (<https://github.com/allenai/scibert>), which lists no dataset-specific licenses but is released with an Apache-2.0 license.

**Paper Field** This problem asks models to predict a paper’s area of study given its title. This task’s dataset contains 84,000/5,599/22,399 train/val/test sentences. Though the original dataset is derived from the MAG [75], it was formulated into this task format by SciBERT directly [91].

**SciCite** This problem tasks models to predict an “intent” label for sentences that cite other scientific works within academic articles. This task’s dataset contains 7,320/916/1,861 train/val/test sentences, and is originally sourced from [104].

**ACL-ARC** This problem tasks models to predict an “intent” label for sentences that cite other scientific works within academic articles. This task’s dataset contains 1,688/114/139 train/val/test sentences and is originally sourced from [105].

### **Further Details on the NETWORKS Dataset and FT tasks**

**PT Dataset** We use a dataset of  $\sim 70$ K protein-protein interaction (PPI) ego-networks here, sourced from [43]. Each individual sample here describes a single protein, realized as a biological network (i.e., an attributed graph) corresponding to the ego-network about that protein (i.e., a small subgraph containing all nodes within the target protein) in a broader PPI graph. Unlike our other domains, this domain does not contain sequences. The NETWORKS PT dataset releases its code and dataset files under an MIT license.

**PT Graph** The dataset from [43] is labeled with the presence or absence of any of 4000 protein gene ontology terms associated with the central protein in each PPI ego network. Leveraging these labels, two PPI ego-networks are linked in  $G_{PT}$  if and only if the Hamming distance between their observed label vectors is no more than 9. This is an alternate-representation nearest-neighbor graph.

**FT Dataset/Tasks** Our FT task is the multi-label binary classification of the 40 gene-ontology term annotations (metric: macro-AUROC) used in [43]. We use the PT set for FT training and evaluate the model on a held-out random 10% split.

**Baselines** We compare against both attribute-masking [43] and multi-task supervised PT.

The Networks FT task is a multi-task, binary classification task. Recall that the dataset here consists of PPI ego-networks, which means that an individual sample input to the model is an attributed graph  $x$  which contains a central node, corresponding to a protein, along with the ego-graph surrounding that node in a larger PPI graph. This ego-graph can thus be seen to correspond to the central protein, and the FT and PT tasks leverage this association, as both of which flag whether or not that central protein is associated with particular gene-ontology (GO) terms (annotations relating to protein properties or function applied in the literature). The PT tasks contain 4000 possible GO annotations, but the FT tasks correspond to a smaller set of only 40 GO terms, chosen as they were of greater interest than the full set. See the original source ([43]) for more information and full details.

### **Further Details on Experimental Procedure**

To minimize computational burden, we do not pre-train a structure-inducing model from scratch for PROTEINS and ABSTRACTS datasets. Instead, we initialize a model from the per-token baseline directly, then perform additional pre-training for only a small number of epochs under the new SIPT loss subdivision. We assess both multi-similarity and contrastive  $\mathcal{L}_{SI}$  variants in these domains. On the NETWORKS dataset, we pre-train all models (including baselines) from scratch, and based on early experimental results, we only assess the contrastive loss variant.

### **Further Details on Ablation Studies**

Note that the warm-start procedure described above on the PROTEINS and ABSTRACTS domains allows a powerful ablation study: by additionally training a PT model from the per-token baseline with  $\lambda_{SI} = 0$ , we can uniquely assess the impact of the new loss term, rather than simply additional training or the different PT dataset. We perform this ablation study for all applicable datasets. For the NETWORKS dataset, no additional ablation studies are needed to assess the impact of the loss term, given all models are trained from scratch with the same early-stop procedures.

### **Further Details on Choosing $\lambda_{SI}$**

For the PROTEINS and ABSTRACTS dataset, to choose the optimal value of  $\lambda_{SI}$  for use at PT time, we pre-trained several models and evaluated their efficacy in a link retrieval task on  $G_{PT} = (V, E)$ . In particular, we score a node embedder  $f$  by embedding all nodes  $n \in V$  as  $f(n)$ , then rank all other nodes  $n'$  by the euclidean distance between  $f(n)$  and  $f(n')$ , and assess this ranked list via IR metrics including label ranking average precision (LRAP), normalized discounted cumulative gain (nDCG), average precision (AP), and mean reciprocal rank (MRR), where a node  $n'$  is deemed to be a “successful” retrieval for  $n$  if  $(n, n') \in E$ . In this way, note that we choose  $\lambda_{SI}$  in a manner that is independent of the fine-tuning task and can be determined solely based on the PT data. Final results for these experiments are shown in Methods Table 9 for the proteins dataset and Methods Table 10 for scientific articles.

Ultimately, this process suggests that  $\lambda_{SI}$  of 0.1 is a robust setting, and as such, 0.1 was used directly for the NETWORKS task without further optimization.

### **Further Details on Architecture & Hyperparameters**

The architectures of our encoders for the PROTEINS and ABSTRACTS domains are fully determined from our source models in TAPE [15] and SciBERT [91]. In particular, for proteins and scientific articles, we use a 12-layer Transformer with a hidden size of 768, an intermediate size of 3072, and 12 attention heads. Provided TAPE and SciBERT tokenizers are also used. A single linear layer to the output dimensionality of each task is used as the prediction head, taking as input the output of the final layer’s [CLS] token as a whole-sequence embedding. We also tested either pre-training for a single or for four additional epochs, based on validation set performance, and ultimately used a single epoch for proteins and four for scientific articles.

For the NETWORKS domain, we match the architecture used in the original source [43] for the mask model runs. Save that for computational efficiency, we scale the batch size up as high as it can go, then proportionally scale up the learning rate to account for the larger batch size. This corresponds to a batch size of 1024, the learning rate of 0.01, a GCNN encoder type of GIN, embedding dimensions of 300, 5 layers, 10% dropout, mean pooling, and a JK strategy of “last”.

Fine-tuning hyperparameters (learning rate, batch size, and the number of epochs) were determined based on a combination of existing results, hyperparameter tuning, and machine limitations. On proteins, most hyperparameters were set to follow those reported for a LM PT model in [106], though additional limited hyperparameter searches were performed to validate that these choices were adequate. As the original source for these hyperparameters was an LM PT model, any bias here should be *against* SIPT, meaning this is a conservative choice. Early stopping (based on the number of epochs without observing improvement in the validation set performance) was employed, and batch size was set as large as possible given the limitations of the underlying machine. For the PLUS reproduction, we compared hyperparameters analogous to the reported PLUS hyperparameters for other tasks and analogous to our hyperparameters for other tasks and used those that performed best on the validation set. For scientific articles, we performed a grid search to optimize downstream task performance on the validation set, with the learning rate varying between 5e-6 and 5e-5 and the number of epochs between 2 and 5. The same grid search was used in the original SciBERT method. We additionally match the SciBERT benchmark by applying a dropout of 0.1, using the Adam optimizer with linear warm-up and decay, a batch size of 32, and no early stopping. For the NETWORKS, FT hyperparameters were again chosen to match the original source model [43] to save the increase in batch size and learning rate. No additional hyperparameter search was performed.

Final hyperparameters for each downstream task are shown in Tables 5 for proteins and 6 for scientific articles.

Task	Batch Size	LR
Remote Homology	16	1e-5
Fluorescence	128	5e-5
Stability	512	1e-4
Secondary Structure	16	1e-5

**Table 5:** Final hyperparameters for our PROTEINS domain. All tasks used 200 total epochs and performed early stopping after 25 epochs of no validation set improvement. LR, learning rate.

Task	# Epochs	LR
Paper Field	2	5e-5
ACL-ARC	4/5	5e-5
SciCite	3/2	1e-5

**Table 6:** Final hyperparameters for our ABSTRACTS dataset. All models used a batch size of 32 and no early stopping to match the original SciBERT paper [91]. LR, learning rate. A / B = [LM PT Hyperparameter] / [SIPT Hyperparameter].

### Further Details on Implementation and Compute Environment

We leverage PyTorch for our codebase. FT Experiments and NETWORKS PT were run over various ubuntu machines (versions ranged from 16.04 to 20.04) with a variety of NVIDIA GPUs. PROTEINS and ABSTRACTS PT runs were performed on a Power 9 system, each run using 4 NVIDIA 32 GB V100 GPUs with InfiniBand at half precision.

### Full Results

Here we provide the raw FT results for all tasks in the PROTEINS and ABSTRACTS domains, respectively (Tables 7 and 8). The NETWORKS domain raw results are already present in the main text (Figure 3).

Model	RH	FL	ST	SS	CP
TAPE	21%	<b>0.68</b>	0.73	73%	0.32
PLUS	19.8%±1.7*	0.63	0.76	73%	N/A
LM PT	23.8%±1.1	0.67±0.00	0.76±0.02	73.9%±0.0	0.38
SIPT-C	25.1%±0.6	<b>0.68±0.00</b>	<b>0.77±0.01</b>	73.9%±0.0	0.38
SIPT-M	<b>26.6%±1.0</b>	<b>0.68±0.00</b>	0.76±0.01	<b>74.2%±0.1</b>	<b>0.39</b>

**Table 7:** Results of the TAPE Transformer [15], the PLUS Transformer [45] (\*: our measurements), our LM PT baseline, and two SIPT variants (“-C” indicates the contrastive loss, “-M” the multisimilarity loss). Higher is better.

Model	PF	SC	AA	SRE
SciBERT	<b>0.66</b>	0.85	0.71	0.80
LM PT	<b>0.66<math>\pm</math>0.0</b>	0.85 $\pm$ 0.01	0.70 $\pm$ 0.05	0.80 $\pm$ 0.01
SIPT-C	<b>0.66<math>\pm</math>0.0</b>	<b>0.86<math>\pm</math>0.01</b>	<b>0.76<math>\pm</math>0.02</b>	<b>0.81<math>\pm</math>0.00</b>
SIPT-M	<b>0.66<math>\pm</math>0.0</b>	0.85 $\pm$ 0.00	0.73 $\pm$ 0.05	N/A

**Table 8:** Results of the original SciBERT [91] model, our own LM PT baseline, and two SIPT variants (“-C” indicates the contrastive loss, “-M” the multisimilarity loss). Higher is better.

### SIPT Results are in Accordance with Theory and Guiding Hypothesis

Results over all real-world domains are consistent with our theoretical analyses and guiding hypothesis. We can also analyze the extent to which induced structure helps non-NLP domains by examining the results of our  $\lambda_{SI}$  tuning procedure. In particular, we find that far less structure-inducing is necessary on our ABSTRACTS dataset ( $\lambda_{SI} = 0.01$ ) than on our PROTEINS dataset ( $\lambda_{SI} = 0.1$ ). This agrees with our guiding hypothesis that per-sample latent space regularization is much more necessary on non-NLP domains than on NLP domains.

To demonstrate this, we show the final results for the guiding link-retrieval task for the PROTEINS domain in Table 9 and for the ABSTRACTS domain in Table 10. In both settings, we compare the following models.

**Random** Nodes are embedded with random vectors to assess chance performance.

**Initial Model** Nodes are embedded with the base pre-trained model we build on in our experiments without further modifications. This model is TAPE [15] for proteins and SciBERT [91] for scientific articles.

**LM PT** Nodes are embedded with the final encoder after additional pre-training on our graph-augmented datasets, but without any SIPT (i.e.,  $\lambda_{SI} = 0$ ).

**CS RoBERTa** (*for scientific articles only*) Nodes are embedded via [12]’s DAPT CS RoBERTa model, which is another LM PT model over scientific abstracts which performed very well on ACL-ARC, the task on which SIPT does best in scientific articles.

**SIPT** (*for various values of  $\lambda_{SI}$* ). Nodes are represented via SIPT PT models at the specified weighting. For proteins, all SIPT models are initialized from TAPE, but for scientific articles, we test against both initializing from SciBERT and CS RoBERTa (as both are just different, domain-specific LM PT models).

Note that in addition to the discrepancy in the magnitude of improvement (over scientific articles, average precision goes from 12.9% to 14.2%, vs. 2.4% to 3.5% on proteins, which is proportionally much more significant), we can also see that SIPT improves retrieval performance

Method	$\lambda_{SI}$	LRAP	nDCG	AP	MRR
Random Baseline	N/A	0.88%	27.1%	0.88%	0.003
TAPE [15]	N/A	8.50%	34.9%	2.41%	0.226
LM PT Baseline	0	8.92%	38.0%	2.33%	0.238
SIPT (TAPE Initialized)	0.01	9.69%	39.1%	2.56%	0.254
	0.10	10.95%	39.4%	3.46%	0.260
	0.50	10.54%	40.3%	3.43%	0.246
	0.90	10.12%	39.0%	3.16%	0.237
	0.99	14.50%	37.5%	3.13%	0.236

**Table 9:** PT set link-retrieval performance for a random baseline, the raw TAPE model, and SIPT for various weighting parameters  $\lambda_{SI}$  on the dataset of protein sequences. LRAP, label ranking average precision; nDCG, normalized discounted cumulative gain; AP, average precision; MRR, mean reciprocal rank. Higher values indicate better performance. Highlighted in grey are realizations of SIPT framework that yield better results than the strongest baseline, providing evidence that incorporating sequence-level relational information into PT (*i.e.*,  $\lambda_{SI} > 0$ ) leads to improved performance.

over the baselines for proteins much more than it does for scientific articles. This is, admittedly, largely due to [12]’s CS RoBERTa model’s surprisingly good performance without any modifications, however as we also compare SIPT pre-trained from a CS RoBERTa model and it does not demonstrate significant improvements, we still feel this is a fair comparison. These findings are consistent with our hypothesis that SIPT will offer more significant advantages in non-natural language domains.

Method	$\lambda_{SI}$	LRAP	nDCG	AP	MRR
Random Baseline	N/A	0.89%	26.0%	0.27%	0.016
SciBERT [91]	N/A	17.22%	52.8%	5.16%	0.272
LM PT Baseline (SciBERT initialized)	0	16.79%	35.4%	5.00%	0.271
DAPT CS RoBERTa [12]	N/A	32.56%	50.3%	12.86%	0.459
LM PT Baseline (CS RoBERTa initialized)	0	30.58%	48.3%	12.36%	0.438
SIPT (SciBERT initialized)	0.01	42.26%	58.7%	14.23%	0.536
	0.10	34.73%	52.5%	9.39%	0.457
	0.50	32.85%	50.8%	8.37%	0.438
	0.90	31.61%	49.8%	7.82%	0.426
	0.99	30.72%	49.0%	6.80%	0.415
SIPT (CS RoBERTa initialized)	0.01	33.32%	51.2%	8.61%	0.448
	0.10	25.46%	44.4%	5.88%	0.359
	0.50	25.08%	44.0%	6.08%	0.355
	0.90	22.43%	41.6%	4.27%	0.317
	0.99	22.38%	41.5%	4.68%	0.316

**Table 10:** PT set link-retrieval performance for a random baseline, the raw SciBERT model, and SIPT for various weighting parameters  $\lambda_{SI}$  on the scientific articles dataset. LRAP, label ranking average precision; nDCG, normalized discounted cumulative gain; AP, average precision; MRR, mean reciprocal rank. Higher values indicate better performance. Highlighted in grey are realizations of SIPT framework that yield better results than the strongest baseline, providing evidence that incorporating sequence-level relational information into PT (*i.e.*,  $\lambda_{SI} > 0$ ) leads to improved performance.

## References

1. Peters, M. E. *et al.* Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2227–2237 (Association for Computational Linguistics, New Orleans, Louisiana, 2018).
2. Brown, T. B. *et al.* Language models are few-shot learners. *arXiv preprint arXiv:2005.14165* (2020).
3. Raffel, C. *et al.* Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research* **21**, 1–67 (2020).
4. Liu, Y. *et al.* RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv:1907.11692 [cs]* (2019). ArXiv: 1907.11692.
5. Radford, A., Narasimhan, K., Salimans, T. & Sutskever, I. Improving language understanding by generative pre-training (2018).
6. Radford, A. *et al.* Language models are unsupervised multitask learners. *OpenAI blog* **1**, 9 (2019).
7. Lewis, M. *et al.* BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 7871–7880 (Association for Computational Linguistics, Online, 2020).
8. Conneau, A. *et al.* Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 8440–8451 (Association for Computational Linguistics, Online, 2020).
9. Clark, K., Luong, M.-T., Le, Q. V. & Manning, C. D. ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators. In *ICLR* (2019).
10. Joshi, M. *et al.* SpanBERT: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics* **8**, 64–77 (2020).
11. Dong, L. *et al.* Unified language model pre-training for natural language understanding and generation. In Wallach, H. *et al.* (eds.) *Advances in Neural Information Processing Systems*, vol. 32 (Curran Associates, Inc., 2019).
12. Gururangan, S. *et al.* Don’t stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 8342–8360 (Association for Computational Linguistics, Online, 2020).
13. Sun, Y. *et al.* ERNIE: Enhanced Representation through Knowledge Integration. *arXiv:1904.09223 [cs]* (2019). ArXiv: 1904.09223.
14. Peters, M. E. *et al.* Knowledge enhanced contextual word representations. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 43–54 (2019).
15. Rao, R. *et al.* Evaluating Protein Transfer Learning with TAPE. In *NeurIPS* (Curran Associates, Inc., 2019).

16. Yamada, I., Asai, A., Shindo, H., Takeda, H. & Matsumoto, Y. Luke: Deep contextualized entity representations with entity-aware self-attention. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 6442–6454 (2020).
17. Sanh, V. *et al.* Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations* (2022).
18. Xiong, W., Du, J., Wang, W. Y. & Stoyanov, V. Pretrained encyclopedia: Weakly supervised knowledge-pretrained language model. In *International Conference on Learning Representations* (2020).
19. Rao, R. M. *et al.* Msa transformer. In Meila, M. & Zhang, T. (eds.) *Proceedings of the 38th International Conference on Machine Learning*, vol. 139 of *Proceedings of Machine Learning Research*, 8844–8856 (PMLR, 2021).
20. Sun, T. *et al.* CoLAKE: Contextualized language and knowledge embedding. In *Proceedings of the 28th International Conference on Computational Linguistics*, 3660–3670 (International Committee on Computational Linguistics, Barcelona, Spain (Online), 2020).
21. He, B. *et al.* BERT-MK: Integrating graph contextualized knowledge into pre-trained language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2281–2290 (Association for Computational Linguistics, Online, 2020).
22. Qin, Y. *et al.* ERICA: Improving entity and relation understanding for pre-trained language models via contrastive learning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 3350–3363 (Association for Computational Linguistics, Online, 2021).
23. Yu, D., Zhu, C., Yang, Y. & Zeng, M. Jacket: Joint pre-training of knowledge graph and language understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, 11630–11638 (2022).
24. Zhou, W., Lee, D.-H., Selvam, R. K., Lee, S. & Ren, X. Pre-training text-to-text transformers for concept-centric common sense. In *International Conference on Learning Representations* (2021).
25. Yuan, Z., Liu, Y., Tan, C., Huang, S. & Huang, F. Improving biomedical pretrained language models with knowledge. In *Proceedings of the 20th Workshop on Biomedical Language Processing*, 180–190 (Association for Computational Linguistics, Online, 2021).
26. Zhang, X.-C. *et al.* Mg-bert: leveraging unsupervised atomic representation learning for molecular property prediction. *Briefings in Bioinformatics* (2021).
27. Caciularu, A. *et al.* CDLM: Cross-document language modeling. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, 2648–2662 (Association for Computational Linguistics, Punta Cana, Dominican Republic, 2021).
28. He, B., Jiang, X., Xiao, J. & Liu, Q. Kgplm: Knowledge-guided language model pre-training via generative and discriminative learning. *arXiv preprint arXiv:2012.03551* (2020).
29. Levine, Y. *et al.* The inductive bias of in-context learning: Rethinking pretraining example design. In *International Conference on Learning Representations* (2022).

30. Li, D., Yi, M. & He, Y. LP-BERT: multi-task pre-training knowledge graph BERT for link prediction. *CoRR* **abs/2201.04843** (2022). [2201.04843](#).
31. BehnamGhader, P., Zakerinia, H. & Baghshah, M. S. Mg-bert: Multi-graph augmented bert for masked language modeling. In *Proceedings of the Fifteenth Workshop on Graph-Based Methods for Natural Language Processing (TextGraphs-15)*, 125–131 (2021).
32. Sun, K., Li, Z. & Zhao, H. Multilingual pre-training with universal dependency learning. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P. & Vaughan, J. W. (eds.) *Advances in Neural Information Processing Systems*, vol. 34, 8444–8456 (Curran Associates, Inc., 2021).
33. Rives, A. *et al.* Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences* **118**, e2016239118 (2021).
34. Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi, M. & Church, G. M. Unified rational protein engineering with sequence-based deep representation learning. *Nature Methods* (2019).
35. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL* (2019).
36. Zhang, Z. *et al.* ERNIE: Enhanced Language Representation with Informative Entities. In *Proceedings of the 57th Annual Meeting of the ACL* (ACL, 2019).
37. Su, Y. *et al.* Cokebert: Contextual knowledge selection and embedding towards enhanced pre-trained language models. *AI Open* **2**, 127–134 (2021).
38. Zhang, Z. & Zhao, H. Structural pre-training for dialogue comprehension. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 5134–5145 (Association for Computational Linguistics, Online, 2021).
39. Kuncoro, A. *et al.* Syntactic structure distillation pretraining for bidirectional encoders. *Transactions of the Association for Computational Linguistics* **8**, 776–794 (2020).
40. Lan, Z. *et al.* ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. In *ICLR* (2019).
41. Zhang, T. *et al.* Smedbert: A knowledge-enhanced pre-trained language model with structured semantics for medical text mining. *arXiv preprint arXiv:2108.08983* (2021).
42. Liu, X., He, P., Chen, W. & Gao, J. Multi-Task Deep Neural Networks for Natural Language Understanding. In *ACL* (2019).
43. Hu, W. *et al.* Strategies for Pre-training Graph Neural Networks. In *ICLR* (2020).
44. Ke, P., Ji, H., Liu, S., Zhu, X. & Huang, M. SentiLARE: Sentiment-aware language representation learning with linguistic knowledge. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 6975–6988 (Association for Computational Linguistics, Online, 2020).
45. Min, S., Park, S., Kim, S., Choi, H.-S. & Yoon, S. Pre-Training of Deep Bidirectional Protein Sequence Representations with Structural Information. *arXiv: 1912.05625* (2020).

46. McDermott, M. B. A. *et al.* A Comprehensive Evaluation of Multi-task Learning and Multi-task Pre-training on EHR Time-series Data. *arXiv: 2007.10185* (2020).
47. Sun, Y. *et al.* ERNIE 2.0: A Continual Pre-Training Framework for Language Understanding. *Proceedings of the AAAI Conference on Artificial Intelligence* **34** (2020). Number: 05.
48. Sun, Y. *et al.* Ernie 3.0: Large-scale knowledge enhanced pre-training for language understanding and generation. *arXiv preprint arXiv:2107.02137* (2021).
49. Yu, W. *et al.* Dict-BERT: Enhancing language model pre-training with dictionary. In *Findings of the Association for Computational Linguistics: ACL 2022*, 1907–1918 (Association for Computational Linguistics, Dublin, Ireland, 2022).
50. Yasunaga, M., Leskovec, J. & Liang, P. LinkBERT: Pretraining language models with document links. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 8003–8016 (Association for Computational Linguistics, Dublin, Ireland, 2022).
51. Wang, W. *et al.* Structbert: Incorporating language structures into pre-training for deep language understanding. In *International Conference on Learning Representations* (2020).
52. Lewis, M. *et al.* Pre-training via paraphrasing. *Advances in Neural Information Processing Systems* **33**, 18470–18481 (2020).
53. Guu, K., Lee, K., Tung, Z., Pasupat, P. & Chang, M.-W. Realm: Retrieval-augmented language model pre-training. In *Proceedings of the 37th International Conference on Machine Learning, ICML'20 (JMLR.org, 2020)*.
54. You, Y. *et al.* Graph contrastive learning with augmentations. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. & Lin, H. (eds.) *Advances in Neural Information Processing Systems*, vol. 33, 5812–5823 (Curran Associates, Inc., 2020).
55. Qiu, J. *et al.* Gcc: Graph contrastive coding for graph neural network pre-training. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '20*, 1150–1160 (Association for Computing Machinery, New York, NY, USA, 2020).
56. Giorgi, J., Nitski, O., Wang, B. & Bader, G. DeCLUTR: Deep contrastive learning for unsupervised textual representations. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 879–895 (Association for Computational Linguistics, Online, 2021).
57. Wu, Z. *et al.* Clear: Contrastive learning for sentence representation. *arXiv preprint arXiv:2012.15466* (2020).
58. You, Y., Chen, T., Shen, Y. & Wang, Z. Graph contrastive learning automated. In Meila, M. & Zhang, T. (eds.) *Proceedings of the 38th International Conference on Machine Learning*, vol. 139 of *Proceedings of Machine Learning Research*, 12121–12132 (PMLR, 2021).
59. Meng, Y. *et al.* Coco-lm: Correcting and contrasting text sequences for language model pretraining. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P. & Vaughan, J. W. (eds.) *Advances in Neural Information Processing Systems*, vol. 34, 23102–23114 (Curran Associates, Inc., 2021).

60. Kong, L. *et al.* A mutual information maximization perspective of language representation learning. In *International Conference on Learning Representations* (2020).
61. Zhang, S., Hu, Z., Subramonian, A. & Sun, Y. Motif-driven contrastive learning of graph representations. *arXiv preprint arXiv:2012.12533* (2020).
62. Carlsson, F., Gyllensten, A. C., Gogoulou, E., Hellqvist, E. Y. & Sahlgren, M. Semantic re-tuning with contrastive tension. In *International Conference on Learning Representations* (2021).
63. Luo, F., Yang, P., Li, S., Ren, X. & Sun, X. Capt: contrastive pre-training for learning denoised sequence representations. *arXiv preprint arXiv:2010.06351* (2020).
64. Zhang, Z. *et al.* Protein representation learning by geometric structure pretraining. *arXiv preprint arXiv:2203.06125* (2022).
65. Chi, Z. *et al.* InfoXLM: An information-theoretic framework for cross-lingual language model pre-training. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 3576–3588 (Association for Computational Linguistics, Online, 2021).
66. Shen, T. *et al.* Exploiting structured knowledge in text via graph-guided representation learning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 8980–8994 (Association for Computational Linguistics, Online, 2020).
67. Fang, Y. *et al.* Molecular contrastive learning with chemical element knowledge graph. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, 3968–3976 (2022).
68. Wang, X. *et al.* KEPLER: A Unified Model for Knowledge Embedding and Pre-trained Language Representation. *Transactions of the Association for Computational Linguistics* **9**, 176–194 (2021). [https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl\\_a\\_00360/1923927/tacl\\_a\\_00360.pdf](https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00360/1923927/tacl_a_00360.pdf).
69. Fang, Y. *et al.* Knowledge-aware contrastive molecular graph learning. *arXiv preprint arXiv:2103.13047* (2021).
70. Jiang, X., Liang, Y., Chen, W. & Duan, N. Xlm-k: Improving cross-lingual language model pre-training with multilingual knowledge. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, 10840–10848 (2022).
71. Guo, Y. *et al.* Webformer: Pre-training with web pages for information retrieval. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 1502–1512 (2022).
72. Deng, J. *et al.* Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255 (Ieee, 2009).
73. Li, B. *et al.* On the sentence embeddings from pre-trained language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 9119–9130 (Association for Computational Linguistics, Online, 2020).
74. Zitnik, M., Sosič, R., Feldman, M. W. & Leskovec, J. Evolution of resilience in protein interactomes across the tree of life. *Proceedings of the National Academy of Sciences* **116** (2019).

75. Wang, K. *et al.* A Review of Microsoft Academic Services for Science of Science Studies. *Frontiers in Big Data* (2019).
76. Hu, W. *et al.* Open Graph Benchmark: Datasets for Machine Learning on Graphs. *arXiv:2005.00687* (2020).
77. Sanh, V. *et al.* Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207* (2021).
78. Saunshi, N., Plevrakis, O., Arora, S., Khodak, M. & Khandeparkar, H. A Theoretical Analysis of Contrastive Unsupervised Representation Learning. In *ICML* (2019).
79. Ribeiro, D. N. & Forbus, K. Combining analogy with language models for knowledge extraction. In *3rd Conference on Automated Knowledge Base Construction* (2021).
80. Gao, H. & Huang, H. Deep attributed network embedding. In *Twenty-Seventh International Joint Conference on Artificial Intelligence (IJCAI)* (2018).
81. Cui, G., Zhou, J., Yang, C. & Liu, Z. Adaptive graph encoder for attributed graph embedding. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 976–985 (2020).
82. Li, Y., Sha, C., Huang, X. & Zhang, Y. Community detection in attributed graphs: An embedding approach. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32 (2018).
83. Li, M. M., Huang, K. & Zitnik, M. Representation learning for networks in biology and medicine: Advancements, challenges, and opportunities. *arXiv:2104.04883* (2021).
84. Kipf, T. N. & Welling, M. Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings* (OpenReview.net, 2017).
85. Hamilton, W. L., Ying, R. & Leskovec, J. Inductive representation learning on large graphs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 1025–1035 (2017).
86. Vert, J.-P. & Yamanishi, Y. Supervised graph inference. In *NeurIPS* (2004).
87. Shaw, B. & Jebara, T. Structure preserving embedding. In *ICML* (2009).
88. Shaw, B., Huang, B. & Jebara, T. Learning a Distance Metric from a Network. In *NeurIPS* (2011).
89. Hadsell, R., Chopra, S. & LeCun, Y. Dimensionality Reduction by Learning an Invariant Mapping. In *CVPR* (2006).
90. Wang, X., Han, X., Huang, W., Dong, D. & Scott, M. R. Multi-Similarity Loss With General Pair Weighting for Deep Metric Learning. In *CVPR* (2019).
91. Beltagy, I., Lo, K. & Cohan, A. SciBERT: A Pretrained Language Model for Scientific Text. In *EMNLP* (2019).

92. Gao, T., Yao, X. & Chen, D. SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 6894–6910 (Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 2021).
93. Saunshi, N., Malladi, S. & Arora, S. A mathematical exploration of why language models help solve downstream tasks. In *International Conference on Learning Representations* (2021).
94. Yoon, J., Zhang, Y., Jordon, J. & van der Schaar, M. VIME: Extending the Success of Self- and Semi-supervised Learning to Tabular Domain. In *NeurIPS* (2020).
95. Zeng, H., Zhou, H., Srivastava, A., Kannan, R. & Prasanna, V. Graphsaint: Graph sampling based inductive learning method. In *International Conference on Learning Representations* (2020).
96. Zhu, J. *et al.* Generalizing graph neural networks beyond homophily. In *NeurIPS* (2020).
97. Huang, K. & Zitnik, M. Graph meta learning via local subgraphs. In *NeurIPS* (2020).
98. Zhang, D., Yin, J., Zhu, X. & Zhang, C. Homophily, structure, and content augmented network representation learning. In *ICDM* (2016).
99. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830 (2011).
100. Hou, J., Adhikari, B. & Cheng, J. DeepSF: deep convolutional neural network for mapping protein sequences to folds. *Bioinformatics* (2018).
101. Klausen, M. S. *et al.* NetSurfP-2.0: Improved prediction of protein structural features by integrated deep learning. *Proteins* (2019).
102. Rocklin, G. J. *et al.* Global analysis of protein folding using massively parallel design, synthesis, and testing. *Science* **357** (2017).
103. Sarkisyan *et al.* Local fitness landscape of the green fluorescent protein. *Nature* **533** (2016).
104. Cohan, A., Ammar, W., van Zuylen, M. & Cady, F. Structural Scaffolds for Citation Intent Classification in Scientific Publications. In *Proceedings of the 2019 Conference of the NAACL: Human Language Technologies, Volume 1 (Long and Short Papers)* (2019).
105. Jurgens, D., Kumar, S., Hoover, R., McFarland, D. & Jurafsky, D. Measuring the Evolution of a Scientific Field through Citation Frames. *Transactions of the ACL* **6** (2018).
106. McDermott, M., Yap, B., Hsu, H., Jin, D. & Szolovits, P. Adversarial contrastive pre-training for protein sequences. *arXiv preprint arXiv:2102.00466* (2021).
107. Févry, T., Baldini Soares, L., FitzGerald, N., Choi, E. & Kwiatkowski, T. Entities as experts: Sparse memory access with entity supervision. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 4937–4951 (Association for Computational Linguistics, Online, 2020).
108. Mroueh, Y., Poggio, T., Rosasco, L. & Slotine, J.-J. E. Multiclass learning with simplex coding. In *Proceedings of the 25th International Conference on Neural Information Processing Systems-Volume 2*, 2789–2797 (2012).

109. Veličković, P. *et al.* Graph attention networks. In *International Conference on Learning Representations* (2018).
110. Wang, Z., Zhang, J., Feng, J. & Chen, Z. Knowledge graph embedding by translating on hyperplanes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 28 (2014).
111. Sun, Z., Deng, Z.-H., Nie, J.-Y. & Tang, J. Rotate: Knowledge graph embedding by relational rotation in complex space. In *International Conference on Learning Representations* (2018).
112. Liu, Y., Wan, Y., He, L., Peng, H. & Philip, S. Y. Kg-bart: Knowledge graph-augmented bart for generative commonsense reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, 6418–6425 (2021).
113. Yang, J. *et al.* Graphformers: Gnn-nested language models for linked text representation. *arXiv preprint arXiv:2105.02605* (2021).
114. Agarwal, O., Ge, H., Shakeri, S. & Al-Rfou, R. Knowledge graph based synthetic corpus generation for knowledge-enhanced language model pre-training. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 3554–3565 (Association for Computational Linguistics, Online, 2021).
115. Lu, Y., Lu, H., Fu, G. & Liu, Q. Kelm: Knowledge enhanced pre-trained language representations with message passing on hierarchical relational graphs. *arXiv preprint arXiv:2109.04223* (2021).
116. Zhang, N. *et al.* Drop redundant, shrink irrelevant: Selective knowledge injection for language pretraining. In *In IJCAI* (2021).
117. Yao, L., Mao, C. & Luo, Y. Kg-bert: Bert for knowledge graph completion. *arXiv preprint arXiv:1909.03193* (2019).
118. Wang, L., Zhao, W., Wei, Z. & Liu, J. SimKGC: Simple contrastive knowledge graph completion with pre-trained language models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 4281–4294 (Association for Computational Linguistics, Dublin, Ireland, 2022).
119. Wang, Z. *et al.* CLEVE: Contrastive Pre-training for Event Extraction. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 6283–6297 (Association for Computational Linguistics, Online, 2021).
120. Kim, W., Son, B. & Kim, I. Vilt: Vision-and-language transformer without convolution or region supervision. In Meila, M. & Zhang, T. (eds.) *Proceedings of the 38th International Conference on Machine Learning*, vol. 139 of *Proceedings of Machine Learning Research*, 5583–5594 (PMLR, 2021).
121. Li, C. *et al.* StructuralLM: Structural pre-training for form understanding. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 6309–6318 (Association for Computational Linguistics, Online, 2021).
122. Fan, Z. *et al.* A unified continuous learning framework for multi-modal knowledge discovery and pre-training. *arXiv preprint arXiv:2206.05555* (2022).

123. Liu, F., Shareghi, E., Meng, Z., Basaldella, M. & Collier, N. Self-alignment pretraining for biomedical entity representations. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 4228–4238 (Association for Computational Linguistics, Online, 2021).
124. Ma, Z. *et al.* *Pre-Training for Ad-Hoc Retrieval: Hyperlink is Also You Need*, 1212–1221 (Association for Computing Machinery, New York, NY, USA, 2021).
125. Calixto, I., Raganato, A. & Pasini, T. Wikipedia entities as rendezvous across languages: Grounding multilingual language models by predicting Wikipedia hyperlinks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 3651–3661 (Association for Computational Linguistics, Online, 2021).
126. Jiang, X., Lu, Y., Fang, Y. & Shi, C. Contrastive pre-training of gnn on heterogeneous graphs. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management, CIKM '21*, 803–812 (Association for Computing Machinery, New York, NY, USA, 2021).
127. Chen, B. *et al.* Code: Contrastive pre-training with adversarial fine-tuning for zero-shot expert linking. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, 11846–11854 (2022).
128. Meng, Z., Liu, F., Clark, T. H., Shareghi, E. & Collier, N. Mixture-of-partitions: Infusing large biomedical knowledge graphs into bert. *arXiv preprint arXiv:2109.04810* (2021).
129. Liu, W. *et al.* K-bert: Enabling language representation with knowledge graph. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, 2901–2908 (2020).
130. Reimers, N. & Gurevych, I. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 3982–3992 (Association for Computational Linguistics, Hong Kong, China, 2019).
131. Yan, Y. *et al.* ConSERT: A contrastive framework for self-supervised sentence representation transfer. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 5065–5075 (Association for Computational Linguistics, Online, 2021).
132. Chen, T., Kornblith, S., Norouzi, M. & Hinton, G. A simple framework for contrastive learning of visual representations. In III, H. D. & Singh, A. (eds.) *Proceedings of the 37th International Conference on Machine Learning*, vol. 119 of *Proceedings of Machine Learning Research*, 1597–1607 (PMLR, 2020).
133. Zhang, Y., He, R., Liu, Z., Lim, K. H. & Bing, L. An unsupervised sentence embedding method by mutual information maximization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1601–1610 (Association for Computational Linguistics, Online, 2020).
134. Faldu, K., Sheth, A., Kikani, P. & Akabari, H. Ki-bert: Infusing knowledge context for better language and domain understanding. *arXiv preprint arXiv:2104.08145* (2021).

135. Yan, R., Sun, L., Wang, F. & Zhang, X. A general method for transferring explicit knowledge into language model pretraining. *Security and Communication Networks* **2021** (2021).
136. Wang, R. *et al.* K-adapter: Infusing knowledge into pre-trained models with adapters. *arXiv preprint arXiv:2002.01808* (2020).
137. Poerner, N., Waltinger, U. & Schütze, H. E-BERT: Efficient-yet-effective entity embeddings for BERT. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, 803–818 (Association for Computational Linguistics, Online, 2020).
138. Gunel, B., Du, J., Conneau, A. & Stoyanov, V. Supervised contrastive learning for pre-trained language model fine-tuning. In *International Conference on Learning Representations* (2021).
139. Zhang, X. *et al.* GreaseLM: Graph REASONing enhanced language models. In *International Conference on Learning Representations* (2022).
140. Khandelwal, U., Levy, O., Jurafsky, D., Zettlemoyer, L. & Lewis, M. Generalization through memorization: Nearest neighbor language models. In *International Conference on Learning Representations* (2020).
141. Kim, T., Yoo, K. M. & Lee, S.-g. Self-guided contrastive learning for BERT sentence representations. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 2528–2540 (Association for Computational Linguistics, Online, 2021).
142. Su, J., Cao, J., Liu, W. & Ou, Y. Whitening sentence representations for better semantics and faster retrieval. *arXiv preprint arXiv:2103.15316* (2021).
143. Huang, J. *et al.* WhiteningBERT: An easy unsupervised sentence embedding approach. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, 238–244 (Association for Computational Linguistics, Punta Cana, Dominican Republic, 2021).
144. Cohan, A., Feldman, S., Beltagy, I., Downey, D. & Weld, D. SPECTER: Document-level representation learning using citation-informed transformers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2270–2282 (Association for Computational Linguistics, Online, 2020).
145. Yan, X., Jian, F. & Sun, B. Sakg-bert: Enabling language representation with knowledge graphs for chinese sentiment analysis. *IEEE Access* **9**, 101695–101701 (2021).

## A Review of Language Model Pre-training Methods

In this supplementary section, we describe all of the models featured in our review (Figure 1 and Table 4) and highlight key details of their approach.

### A.1 Language modelling alone

- [1] General domain NLP; ELMO leverages a biLSTM to perform language modelling; unlike later methods, for FT tasks, models do not typically re-train the entire LSTM but rather use a weighted combination of model interior hidden states as (at FT time) static word-embeddings.
  
- [4] General domain NLP; RoBERTa includes only a masked language modelling objective.
  
- [2,5,6] General domain NLP; The GPT series of models use autoregressive language modelling alone and focus on generative language tasks, not general PT/FT, though GPT-III does show that by reframing many classical NLP fine-tuning tasks as generative language tasks, GPT-III can still offer a compelling zero and few-shot solution to these tasks using only the pre-trained embedder [2].
  
- [7] General domain NLP; BART utilizes a denoising language-model objective across various noising constraints.
  
- [11] General domain NLP; UniLM integrates several different kinds of language modelling, including bidirectional, unidirectional, and sequence-to-sequence LMs. They impose no other PT losses.
  
- [15,33,34] Protein sequences; Various methods have explored language modelling alone for protein sequences. One notable entry is the TAPE benchmark, which also introduces a public benchmark of FT tasks for future comparisons.
  
- [26] Molecular Graphs; Molecular Graph BERT (MG-BERT; no relation to MG-BERT [31]) uses masked atom prediction to pre-train a GNN over molecular graphs.
  
- [8] General domain NLP; This paper pre-trains a model for multi-lingual language modelling, using only a multi-lingual masked language modelling objective.
  
- [12] General domain NLP; DAPT advocates for continual pre-training on increasingly task-focused text to improve its relevance to various downstream tasks. DAPT uses a RoBERTa baseline pre-training model, which includes only a masked language modelling objective. It shows significant gains after adaptation. However, as they only adapt the pre-training context to the more focused text, this induces no additional constraints on the latent space geometry.

[10] General domain NLP; SpanBERT changes the traditional masked language modelling task to a task in which contiguous spans are masked wholesale, rather than individual tokens.

## A.2 Language modelling & templated tasks/prompting as language modelling

[3] General domain NLP; T5 not only performs a robust analysis of various existing pre-training strategies but also introduces the “text-to-text” style of diverse pre-training, in which various downstream NLP tasks can be re-realized as language modelling tasks via templating and prompting, then integrated into language model pre-training alongside unsupervised objectives (such as traditional masked language modelling, albeit realized as a sequence-to-sequence task). As they realize all these downstream tasks as additional language modelling tasks, they neither officially produce a directly constrained per-sample embedding nor constrain the geometry of  $\mathcal{Z}$  beyond traditional masked language modelling.

[24] General domain NLP; CALM builds on ideas from T5 to propose a text-to-text pre-training objective that leverages recognized per-token KG entities from the source text as a generative prompt.

[17] General domain NLP; T0pp extends the architecture of T5 [3] to ingest templated language modelling task from a wide variety of possible input tasks, then evaluates its performance in a zero-shot manner on unseen fine-tuning tasks.

## A.3 Language modelling & Per-token KG Integration

[13] General domain NLP; ERNIE 1 augments traditional MLM with entity-specific masking (e.g., masking the word “Mozart” from the sentence “Mozart was a musician”) to force the model to recover common-sense knowledge about named entities.

[28] General domain NLP; KgPLM adapts the discriminative training ideas of ELECTRA [9] alongside the idea of entity masking explored previously. They perform entity masking and a discriminative loss identifying which tokens were replaced focused on entity replacements.

[22] General domain NLP; ERICA presents a mechanism for leveraging contrastive learning and distant supervision to incorporate external knowledge into a PLM for improving language understanding. ERICA augments MLM with two per-token tasks to ensure the per-token representations within a document reflect the structure of the KG. First, ERICA ensures that the pooled representations of head and tail entities are similar when conditioned on a relation (which is prepended to the document prior to embedding). Second, ERICA ensures that relation embeddings (defined as concatenated head, tail per-token entity embeddings)

are similar within and across documents. As both tasks are done on per-token embeddings and never at a per-sample level, this approach induces minimal constraints on the per-sample latent space.

- [14] General domain NLP; Know-BERT integrates per-token entity information into an MLM pre-training scheme by performing unconstrained attention over a per-entity knowledge graph (only on pre-identified candidate entity spans), alongside any available entity linking supervision information via direct Named Entity Linking. This has similarities with [25] and [107].
- [25] Biomedical domain NLP; KeBioLM integrates a per-token KG into a biomedical language model by augmenting token entity representations with attention lookups into a biomedical KG (regardless of whether the attended entities match a given entity mention in the source text, though they do only apply this on recognized entities). To ensure this attention is meaningful, they perform named entity linking and recognition as auxiliary PT objectives, leveraging the same KG embeddings used during the attention calculation. In doing so, the method incentivizes per-token representations to be similar to their associated entity representations, thus ensuring that the entities are reflected in the attention over the KG. KG embeddings are initialized using Trans-E [108]. Their usage of automatically attending over entities within their language model (without explicit constraints on those matches) is motivated by [107]’s work in [107] and has similarities to Know-BERT [14].
- [16] General domain NLP; LUKE performs pre-training using MLM and an entity-specific masking/recognition scheme that is a slight variation on the traditional entity-specific masking [13] proposed. At FT time, they have other knowledge-specific integrations, including specialized query matrices in KQV attention based on attending to either traditional tokens or entities. However, at PT time, LUKE’s only modulation over a ROBERTA [4] baseline is an entity masking task.
- [20] General domain NLP; COLAKE performs a priori entity linking on the source text, then replaces per-token mentions with entity embeddings, and appends to the input text sub-graphs from a (relational) knowledge graph, including both neighboring mentions and relations in the augmented input text block. This input is then encoded via a transformer that limits attention flow between tokens of different types and trains the entire ensemble with masked language, entity, and relation modelling.
- [18] General domain NLP; In this paper, traditional masked language modelling is augmented with an entity-replacement-detection task. Named entity recognition and linking are performed before pre-training, and entity replacements are constrained to be the same type as

the true entity.

- [30] Knowledge Graph Completion; LP-BERT constructs a specialized pre-training corpus consisting of entity-relation statements from a knowledge graph. This is used in a pre-training context under three pre-training tasks: masked language modelling, masked entity modelling, and masked relationship modelling. All three are per-token, and no per-sample tasks are used at pre-training time.
- [32] Multilingual Language Models; UD-PrLM examines multilingual pre-training, and aims to improve it by incorporating universal dependency parse trees into the model. They incorporate a per-token task to align tokens with identified dependency parse tree components, alongside masked language modelling.

#### **A.4 Language modelling, Per-token KG Integration, & Supervised Classification**

- [47,48] General domain NLP; ERNIE 2.0 & 3.0 augments traditional MLM with entity-specific masking (e.g., masking the word “Mozart” from the sentence “Mozart was a musician”) as well as a multi-task per-sample task, largely motivated at classifying a block of text based on internal text cohesion (predict the true order of the sentences within an input sample & identify whether the sentences within the input sample are spatial neighbors, come from the same document, or come from different documents). ERNIE 3.0 additionally augments pre-training with a per-token relation-embedding task using cloze-filling as a vehicle to perform relation extraction on pre-specified per-token KGs.
- [36] General domain NLP; ERNIE (no relation to [13,47]) uses both architectural and objective-function changes to inject per-token knowledge into PT. Specifically, they separately embed all named entities in a sample using the architecture to join contextualized entity embeddings alongside the embeddings of tokens, realizing that entity in the span and performing entity-specific masking. In addition, they simultaneously perform standard MLM and next-sentence prediction in the manner of BERT [35].
- [31] General domain NLP; MG-BERT introduces a GCNN layer after BERT token, aggregating token embeddings together over a unified graph consisting both of co-occurrence relationships and knowledge graph relationships.
- [23] General domain NLP; JAKET embeds entities by extracting per-token representations of entity texts inside per-entity descriptions, then produces updated KG embeddings via a graph attention network [109]. Those embeddings are then fed into a language model alongside per-token embeddings corresponding to those entities. The entire model is trained according

to an MLM objective, plus entity category prediction and relation prediction (only on the entity embeddings extracted from entity descriptions and fed through the GCNN—not on the raw entities within the contextualized text).

- [21] Biomedical NLP; BERT-MK introduces a transformer-based subgraph summarization network that produces entity embeddings for dynamically chosen subgraphs of a given knowledge graph. This network is trained via a contrastive triplet-validity objective. These are then fused with per-token embeddings in free-text based on apriori entity-token matching (*i.e.*, named entity recognition and linking must be performed first and separately before using this model).
- [37] General domain NLP; Coke is similar to ERNIE [36], JAKET [23], and BERT-MK [21] in that it aggregates entity information by leveraging a GCNN over a restricted dynamic context KG based on token-entity mentions then integrates those augmented embeddings into the per-token embeddings of a BERT-style pretrained model (similar to JAKET and BERT-MK), but also leverages the denoising entity autoencoder task of ERNIE [36]. In addition, in the variant of COKE derived from the BERT model, COKE also employs the next-sentence prediction task introduced in BERT [35].
- [41] Medical domain NLP; SMedBERT leverages a complex, multi-faceted loss including MLM, Sentence-order prediction SOP (as introduced in, e.g., ALBERT [40]), and includes per-token KG information by aggregating token embeddings across KG embeddings (produced via trans-H [110]) corresponding to matching entities and the neighbors of matching entities in the KG. They also include relation and entity masking variations to ensure the PT model learns per-token information corresponding to the KG. This method bares similarity to Coke [37] and JAKET [23]. However, unlike Coke and JAKET, SMedBERT realizes the entity/neighbor matching via a geometric objective, which results in an explicit per-token knowledge graph alignment.
- [49] General domain NLP; Dict-BERT focuses on augmenting BERT by concatenating definitions of rare words via a per-token KG integration. They add two additional tasks atop the traditional MLM task. First, a task maximizing the mutual information between a masked rare word (treated as a named entity) and its definition (represented as the per-token embedding of the first mention of the entity in the concatenated definition). Second, a task discriminating valid rare word definition per-sequence embeddings from non rare-word definition embeddings via a classification objective.
- [44] Sentiment Analysis; SentiLARE integrates sentiment analysis and labels into pre-training by including word polarity signals during masked language modelling and embedding and

augmenting pre-training with a supervised sentence sentiment prediction. Word polarities are determined via an external knowledge base integrated at the per-token level.

- [38] Dialogue Modelling; SPIDER augments traditional MLM and NSP pre-training with two tasks specific to dialogue modelling: first, utterance order prediction, in which individual utterances (which are nested within a larger sample) are shuffled and the true order is predicted, and a geometric task ensuring that subject, verb, object triples from the utterances obey a geometric relationship inspired by KG embedding methods.

## A.5 Language modelling & Graph link-prediction realized as single-task classification

These methods all employ some variant of a graph link-prediction task over their data. However, they all realize this link prediction task not by enforcing any relationship between independent sample embeddings but rather by concatenating samples corresponding to linked (or unlinked, for negative samples) pairs of vertices in the source graph, then framing the learning problem as a binary or multi-class classification problem over the (now concatenated) single output whole sample embedding. In doing so, they transform the task from one that implies a deep geometric constraint over the output latent space to one that only enforces an intra-sample objective and imposes only a shallow geometric constraint on the per-sample latent space.

- [35] General domain NLP; Masked language model plus the binary classification of whether the input text block is sequentially consistent, with samples chosen via true positive pairs vs. randomly joined sentences. This can be seen as a link prediction task over a graph consisting of independent, disconnected “sticks”, with each stick corresponding to sentences in the documents in the corpus, in sequential order.
- [40] General domain NLP; Masked language model plus the binary classification of whether the input text block is sequentially consistent, with samples chosen via true positive pairs vs. reordered positive sentence pairs. This can be seen as a link prediction task over a directed graph consisting of independent, disconnected “sticks”, with each stick corresponding to sentences in the documents in the corpus, in sequential order, with edge direction indicating sequential ordering.
- [50] General domain NLP; Masked language model plus the classification of whether the input text block contains sentences from either (1) random documents, (2) a sequentially consistent pair within a single document, or (3) within a pair of sentences within two linked documents according to a document linking graph  $G$ . This can be seen as a link prediction/edge classification task over a graph whose nodes are text blocks in the corpus, with two distinct

edge modalities. First, to capture sequential consistency within a document, one edge type produces a set of independent, disconnected “sticks”, with each stick corresponding to sentences in the documents in the corpus, in sequential order. Second, to capture the document linking graph  $G$ , sentences in a document  $D_i$  are all linked to all sentences in a document  $D_j$  if and only if documents  $i$  and  $j$  are linked in  $G$ .

- [39] General domain NLP; While this model incorporates an interesting per-token syntactic knowledge distillation procedure, at a per-token level it merely leverages BERT’s NSP loss [35].

## A.6 Language modelling & Single-task Classification

- [45] Protein sequences; Masked language model plus the multi-class classification of to which protein family an input sequence belongs. Uses non-standard whole-sequence embedding procedure (no [CLS] token).
- [51] General domain NLP; StructBERT includes masked language modelling, a token permutation language modelling task, and an extended version of the NSP/SOP task at a per-sample level.

## A.7 Language modelling & Multi-task Classification

- [42] General domain NLP; Masked language model plus multi-task classification across various NLP tasks.
- [43] Graph data; This model uses a masked imputation task similar to a masked language model and a highly multi-task supervised whole-graph level prediction. On this non-NLP domain, [43] finds that the multi-task whole-graph level task is essential for performance.
- [46] EHR Timeseries data; This model uses a masked imputation task similar to a masked language model over time series data and a multi-task supervised whole-sequence prediction task. On this non-NLP domain, [46] finds the multi-task whole-sequence level task essential for performance.

## A.8 Language modelling & whole-sample graph-based contrastive objectives

- [68] General domain NLP; KEPLER augments traditional MLM on text samples with a constraint ensuring the (per-sample) embeddings of entity descriptions pulled from pre-specified knowledge graphs (KGs) reflect geometric constraints, leveraging the [111] geometric constraints. As we will see in our theoretical analyses, these constraints are much more restrictive on the latent space geometry and thus imply a greater encoding of domain knowledge

in the model. Note that JAKET [23] also leverages entity descriptions in its per-token encoding. However, these descriptions are (1) extracted via per-token embeddings, using the first mention of the token, not whole-sample embeddings, and (2) integrated back into the original text in a per-token manner, not optimized over directly via geometric constraints as in KEPLER.

[69] Molecules; CK-GNN designs a pre-training scheme for molecular graphs in which a molecular GNN is trained to produce molecule embeddings that obey the similarity structure of a 1-NN graph in a cluster-limited molecular fingerprint space (using the Dice similarity coefficient). Unlike the NLP approaches, this method has no intra-sample (*i.e.*, per-token, where here “token” refers to individual atoms within the molecular graph) pre-training task.

[70] Multi-lingual NLP; Much like KEPLER, XLM-K augments traditional MLM with two tasks that constrain the geometry of the per-sample latent space via a (now multi-lingual) graph of entity descriptions linked to sentences containing said entities. Like KEPLER, as the graph connections here are defined only for entity descriptions and not all free-text, the latent space regularization is only over a limited slice of the space.

[71] General domain NLP/IR; WebFormer designs a pre-training scheme leveraging the structure of DOM trees in HTML pages to impose multiple per-sample and per-sample/per-token hybrid constraints that encourage individual samples to be (a) close to noised versions of themselves based on reordering or masking and (b) to be close to representations of their parent/child nodes in the DOM tree, thus imposing a structural penalty geometrically. By mixing per-sample and per-token tasks, WebFormer even more closely entangles the per-sample and per-token latent spaces in their model, and this approach bears closer study in other contexts.

## A.9 Language modelling & whole-sample augmentation/noising based contrastive objectives

[60] General domain NLP; InfoWord incorporates an objective alongside masked language modelling which pushes the whole-sample embedding of a sentence to have high mutual information with various sub-contexts within that sentence and low mutual information with sub-contexts of other sentences.

[56] General domain NLP; DeCLUTR optimizes for masked language modelling alongside a contrastive objective comparing anchor spans to positive spans chosen from within individual samples, contrasted against spans from other samples. This is considered “whole-sample” rather than a per-token contrastive loss as the embeddings of the spans (which can be quite long) are produced via a canonicalized pooling operation used for sentence embeddings.

- [57] General domain NLP; CLEAR optimizes for masked language modelling alongside a contrastive objective powered by per-sentence noising strategies, including word or span deletion, reordering, and synonym substitution.
- [59] General domain NLP; COCO-LM builds on other discriminative language modelling variants such as ELECTRA [9] by adding two additional tasks. First, a true language modelling task atop the auxiliary-model-driven corrupted input text. Second, a contrastive objective pushing corrupted sentences towards their un-corrupted originals and those derived from distinct sentences farther apart.
- [62] General domain NLP; Semantic re-tuning via contrastive tension adds a pre-training objective onto language model pre-training. This is done to encourage the final per-sample representations of a single sentence embedded via two otherwise independently trained models to be similar and those of different sentences to be distinct.
- [54,55,58,61,67] Networks; KCL, GraphCL, JOAO, MICRO-Graph and GCC use augmentation-based contrastive learning pre-training methods for network datasets. KCL is notable as it is (1) specialized for molecular graphs and (2) uses a knowledge-derived augmentation strategy that constructs a knowledge enriched version of an input molecular graph as its “augmentation policy.” MICRO-Graph is also notable as its contrastive objective compares a graph to dynamically clustered “motif” subgraphs from within said graph as positive pairs.
- [66] General domain NLP; GLM integrates a per-token KG through traditional entity masking (albeit with an improved selection mechanism) and a per-sample contrastive objective that uses the entity knowledge graph to generate distractor negative samples for the contrastive learning task.
- [63] General domain NLP & Computer Vision; CAPT proposes a noising based contrastive learning loss *in substitution for* the masked language modelling loss of BERT. They employ no per-token pre-training task.
- [64] Protein Sequences/Structures; GearNet introduces a vehicle for pre-training not over protein sequences, but rather over protein structures, realized as graphs. They combine intra-sample/per-amino-acid tasks, including prediction of masked node features and prediction of geometric relationships between nodes as implied by the protein graphs, and a per-sample noising based contrastive objective.

## A.10 Language modelling & multi-modal or multi-lingual contrastive objectives

Note that by viewing multiple data modalities as “augmentations” of the data samples, one can realize these methods (in general) as examples of augmentation-based contrastive learning objectives, such as those used in [92]. However, as these methods are common, we highlight them explicitly here.

- [65] General domain NLP; InfoXLM focuses on multi-lingual pre-training, and leverages per-token tasks. This includes multi-lingual masked language modelling and translation language modelling (*i.e.*, variations on a traditional masked language modelling task). It also incorporates a cross-lingual per-sample contrastive objective that aligns the geometry of the latent spaces across distinct languages. One important nuance is that they use different layer depths to define the latent space for their cross-lingual contrastive objective vs. their per-token objectives, which is not natively describable in our framework. In addition, as each monolingual corpus lacks any rich, independent per-sample task, any individual monolingual latent space cannot be guaranteed to have any rich structural constraints.

## A.11 Language modelling alone with relationally-concatenated samples

These methods concatenate samples together before processing them with a pre-training encoder based on inter-sample relations. This is an orthogonal direction to adding greater per-sample dependencies to pre-training methods than our framework but warrants commentary nonetheless.

- [19] Protein sequences; MSA transformers extend protein-sequence language models such that they do not take in as input a single sequence but rather an entire multiple-sequence alignment (MSA) profile. These profiles consist of many sequences corresponding to evolutionary homologs of the same protein. This concatenated input is processed via a sparsified form of axial self-attention, which enables cross-attention between the various aligned sequences. They impose no per-sequence tasks by default in this architecture.
- [29] General domain NLP; This theoretical analysis shows that transformers cannot model dependencies between sentences that never appear in the same example during pre-training. To combat this, they propose concatenating samples via inter-sample relations (in particular, via a kNN method) at pre-training time, enabling a greater diversity of cross-attention contexts during pre-training vs. fine-tuning. Thus, while they only use language modelling during pre-training, they speculate that their sample-augmentation procedure helps the model better reason about per-sample information through per-token tasks.
- [27] General domain NLP; CDLM proposes to concatenate multiple related documents (leveraging categorical information to cluster documents) together into a single sample prior to

performing traditional masked language modelling. To limit the model’s complexity, attention is restricted to intra-document for unmasked tokens but allowed to be global for masked tokens.

- [53] General domain NLP; REALM uses a latent variable model to learn a relevance score between input text spans and documents in an auxiliary document base. The top- $k$  documents, according to this relevance score, are then concatenated to the input prior to solving the masked language modelling task used during pre-training. In this way, the model learns to join relevant documents from an external knowledge base in accordance with which documents would most improve the masked language modelling objective. In addition, by learning this relevance score, the model introduces an implicit whole-sample structural constraint on the latent space according to the unsupervised clustering induced by relevance assignment.

## A.12 Autoencoding & Unsupervised Clustering

- [52] General domain NLP; MARGE deviates significantly from the norm by not employing any form of language modelling or other forms of a per-token pre-training task. Instead, it employs only a per-sample contextualized autoencoding objective and an unsupervised per-sample retrieval step (to provide context for said autoencoding). While this approach does provide a deeper form of a per-sample structural constraint than many other approaches, it is also implicit and has no mechanism for injecting domain knowledge. MARGE is also tested solely on downstream tasks at the per-sample level, so it is unclear if this method would offer reduced benefits for per-token downstream tasks.

## A.13 Methods orthogonal to our framework

- [112] KG-BART is a text-generation model that leverages per-token knowledge after a text-encoder to enrich the generated text with information from a textual knowledge graph (in a per-token manner). It is neither used for general pre-training nor does it leverage any additional per-sample constraints.
- [113] Text-based Knowledge Graphs; This work produces embeddings of nodes in KGs by combining transformer-based text encodings with graph convolutional network KG embedding methods, leveraging link prediction as the pre-training task. Entity descriptions / textual features represent the individual nodes. Link prediction can be seen as inducing a geometric constraint via the connectivity of the knowledge graph on whole-sample embeddings. However, given that relationships are used in encoding the data as well, GraphFormer cannot be used in a context where KG links may not be observed at FT time. It should be seen not

as a general text PT method but as an advanced KG embedding mechanism, so it does not directly fall under our framework.

- [114] KeLM (unrelated to KELM [115]) is a method for converting a free-text KG into textual nodes so language modelling can be used over that corpus and is orthogonal to the methods of pre-training.
- [79] This paper is a method for populating a KG from free-text via BERT. It has no bearing on incorporating structure or knowledge into PT and is irrelevant to our framework.
- [116] This paper presents a method to drop redundant triples from a knowledge graph and a regularization technique to limit the impact of added irrelevant knowledge to per-token knowledge-enhanced PT methods such as ERNIE [36].
- [117] Knowledge Graph Completion; KG-BERT is a method for knowledge graph completion in which textual representations of entities and relations in KGs are embedded by fine-tuning a pre-trained BERT style transformer for link prediction over a given KG. As this is only for knowledge graph completion, it is orthogonal to our study of pre-trained models in general.
- [118] Knowledge Graph Completion; Much like KG-BERT, SimKGC is a method for knowledge graph completion that fine-tunes a BERT model via a contrastive loss over a fixed knowledge graph for link prediction. Though their methodology overlaps with ours in that both use variants of contrastive losses and SimKGC explores more complex negative sampling strategies, the two methods are still very different. Ours is focused on general pre-training and uses a single encoder and a unified latent space. In contrast, SimKGC is only examined for KG completion and encodes head and tail entities via separate encoders.
- [119] Event Extraction (EE); CLEVE designs a pre-training method specifically for event extraction. Their pre-training method includes a text-encoder which includes a *cross-event* contrastive loss pushing *individual tokens* from the same “event” closer together than those from different events, which bears a surface similarity to our approach. In addition, they add a graph encoder over the semantic structure of events. Their methodology is focused solely on EE, which is orthogonal to our more general PT framework.
- [120] General domain NLP and Computer Vision; ViLT is a method for pre-training aligned text-image pairs. It leverages masked language modelling, an image-text matching binary classification objective, and a contrastive objective comparing image and text representations. This multi-modal contrastive objective is very similar (insofar as it relates to our framework) to those works that perform multi-lingual or other multi-modal contrastive methods. In ViLT, however, the transformer architecture processes images and text jointly in a single encoder,

so it is not well suited for use on only images or only text. This, combined with its focus on computer vision, renders it orthogonal to our framework.

- [121] General domain NLP and Computer Vision; StructuralLM proposes a new method of pre-training for scanned documents that takes advantage of the structure of the document w.r.t. images and text simultaneously. As their focus is on cross-modal pre-training of text and image alignment, it is orthogonal to our work.
- [122] General domain NLP and Computer Vision; This paper proposes a framework for simultaneous (and continuous) discovery of edges in a multi-modal knowledge graph and the leveraging of that knowledge graph to inform representation learning. However, it is not suitable for our framework for two reasons. First, like ViLT, it is focused on image-text alignment pre-training. Second, when producing node (*e.g.*, images or text snippets) representations, it requires connectivity information in the associated multi-modal knowledge graph. In contrast, our methods take as input only elements from  $\mathcal{X}$ .
- [123] Named Entity Linking; SapBERT is a method for aligning the output of a pre-trained language model with a per-token knowledge graph through a metric learning loss applied at a per-sample level but only over entity names (not even entity descriptions). As it applies this as a secondary, post-PT stage, and this method only optimizes for alignment between entity names and a static KG, it is not a general PT framework. It is thus orthogonal to our efforts here.
- [124] Information Retrieval; HARP is a method for specializing pre-training towards ad-hoc query information retrieval. They introduce four retrieval-specific pre-training tasks leveraging hyperlinks in Wikipedia articles in addition to traditional masked language modelling. Rather than using the raw text of the hyperlinks or the per-sample representations of text spans containing hyperlinks, both of which are explored in [125], these authors use attention weights to extract various “queries” from the underlying text and match those against possible destination pages via contrastive losses. This, therefore, does not impose a constraint on the latent space over the original pre-training dataset  $\mathcal{X}$  (but instead introduces a new latent space consisting of query spans) and is further specialized exclusively for ad-hoc retrieval tasks.
- [126] Node Embedding for Heterogeneous Graphs; CPT-HG is a contrastive pre-training framework to embed nodes in a heterogeneous network. Unlike in our setting, where the pre-training graph  $G_{PT}$  is *only used as an implicit input to derive the loss function*, in CPT-HG the graph (with entire edge connectivity information) *is* the input to the problem. Thus,

node embeddings will rely on connectivity information, which is not permissible in our pre-training context. So, this method is orthogonal to our study here.

[127] Expert Matching; CODE is a method specifically and exclusively designed to discover appropriate experts in an employment/contracting setting and is thus orthogonal to our framework, which is focused on more general pre-training.

#### **A.14 Methods that only change things at FT time**

[128] Biomedical domain NLP; MOP does not change anything at PT time but trains sub-KG adapters on entity recognition tasks prior to FT to infuse entity knowledge into the PT method. It is a per-token pre-training method.

[129] General domain NLP; K-BERT, at PT time, is actually equivalent to BERT [35]. However, it does do other interesting things at FT time, including augmenting the sentence flow with injected per-token knowledge graphs and limiting self-attention to only flow along links supported by the original sentence or the injected knowledge. However, as this is only true at FT time, it is equivalent to BERT at PT time.

[130] General domain NLP; This model, at PT time, is equivalent to BERT [35]. Like [129]. However, it specializes in a fine-tuning procedure for sentence information retrieval tasks, similar to how PT is adapted in this framework.

[131] General domain NLP; ConSERT adds an auxiliary specialization stage after pre-training to fine-tune sentence representations. This new stage imposes a SimCLR [132] style data-augmentation/noise-invariance based contrastive learning objective, using adversarial perturbations, token shuffling, token/feature/span erasure, and dropout noising methods.

[133] General domain NLP; IS-BERT does not modify anything from traditional BERT at pre-training time. However, they add a second PT stage to optimize sentence representations alone using an auxiliary feature extractor in the form of various CNNs applied atop BERT token representations. The final sentence representation is trained to maximize mutual information with various sub-contexts within the sentence but low mutual information with other sentences. In this second pre-training stage, there is no language modelling performed. As this approach only adapts an auxiliary featurizer to produce sentence encodings and is not intended for general transfer learning, it is inappropriate for our framework. A similar work that integrates both components during pre-training, and thus is relevant in our work is [60] and is discussed above.

[115] General domain NLP; KELM does not modify PT objective but instead enhances a model at FT time by injecting per-token knowledge via a GNN module atop the pre-trained LM

embeddings via a unified text-entity graph. It is similar to KBERT [129] in this way but resolves other issues with that approach relating to knowledge ambiguity and by supporting multi-hop reasoning, again over the per-token embeddings.

- [134] General domain NLP; KI-BERT augments BERT with KG-specific information via joint token-entity embeddings and information fusion but does this only at FT time.
- [135] General domain NLP; K-XLNet introduces a secondary FT stage in which knowledge injectors throughout an XL-Net architecture are further trained to leverage knowledge (encoded via free-text entity descriptions) that is injected into input sentences alongside matched tokens. It does not modify the XL-Net PT stage at all.
- [136] General domain NLP; K-Adapter proposes to pre-train various knowledge adapters that can be used alongside a pre-trained language model at a fine-tuning time. Thus, while there is a pre-training process for the adapters, this process does not modulate the original pre-trained language model. In addition, both adapters pre-trained in this work are based on per-token knowledge graphs; one leverages concatenated entity embeddings to perform relation classification, and another predicts which token in the sentence is the “head” in a dependency parse tree, so no per-sample constraints are applied.
- [137] General domain NLP; E-BERT injects per-token knowledge into BERT by first aligning embeddings of a knowledge graph with the input word piece embedding space of a (fixed, pre-trained) BERT model, then using various strategies to input them alongside their source mentions in FT text. They do no additional pre-training, so this model only affects the model at FT time.
- [138] General domain NLP; [138] augment LMPT methods with an additional, pre-FT procedure in which the model is further trained using a supervised, per-sample metric learning task leveraging FT labels directly to form the classes used for metric learning. They do not materially change the task-independent PT procedure, though their FT metric learning procedure does induce some structure at the per-sample level.
- [139] QA; GreaseLM is a method for fusing information from knowledge graphs into pre-trained language models. It shares many similarities with methods that do this for pre-training purposes, such as JAKET [23], CokeBERT [37], SMedBERT [41], and Bert-MK [21]. However, unlike these methods, it only employs these techniques at the fine-tuning time, for question answering tasks specifically. As it is not focused on general pre-training, it is outside our scope.

- [140] Language modelling; kNN language models improve the text generation powers of language models by augmenting traditional decoding with a nearest-neighbor lookup operation over a text datastore leveraging the embeddings of a token’s leftward context by the language model to judge nearest neighbors. However, it involves no additional language model training and can only be applied at the fine-tuning time to aid in text generation, and is thus out of our scope.
- [141] Sentence embedding; NT-Xent proposes a secondary specialization stage after pre-training only for generating sentence embeddings. To do this, they employ a contrastive objective contrasting the final CLS embeddings of an updating, specialized BERT model against a pooled aggregate of the per-token embeddings across all layers of the pre-trained BERT model used to initialize the specialized sentence embedding model.
- [73, 142, 143] Sentence Embedding; These methods propose to use unsupervised per-sample smoothing operations (a normalizing flow network in [73] and a mean/covariance standardization whitening operation in [142, 143]) on the per-sample embeddings after pre-training in order to produce higher quality per-sample embeddings.
- [92] General domain NLP; SimCSE extends traditional MLM by imposing a second pre-training stage for optimizing sentence embeddings. In this stage, SimCSE optimizes the transformer such that the whole-sample embeddings satisfy either a supervised or unsupervised contrastive learning objective. In the supervised case, this is based on labeled sentence pairs according to a Natural Language Inference (NLI) task, with entailment pairs being treated as positives and contradiction pairs as hard negatives. In the unsupervised case, this is based solely on applying multiple dropout masks to the same sentence to generate positive pairs. Any two distinct sentence inputs are treated as negative samples. This extra pre-training stage is applied to a relatively small number of samples ( $10^6$ ) relative to the entire PT cohort, which may help prevent catastrophic forgetting of the original pre-training objective.
- [144] Academic NLP; SPECTER extends traditional language model pre-training by imposing a second pre-training stage for optimizing document embeddings (realized as [CLS] token embeddings of concatenated academic paper titles and abstracts). This stage uses a triplet-based geometric loss to ensure that these per-sample embeddings reflect the structure of a pre-specified citation network. This is a form of an explicit, structural constraint; however, they do not ever test fully fine-tuning the SPECTER model in their paper and only compare it against other, frozen pre-trained language models. This is likely to have a significant impact on model comparisons. Similar to SimCSE [92], this extra pre-training stage is applied to a small number of samples (146K documents) to help prevent catastrophic forgetting of the original pre-training objective.

- [125] General domain NLP; This paper introduces a second pre-training stage after multi-lingual masked language modelling. In this second stage, hyperlinks in the source text (drawn from Wikipedia) are matched via single-task classification to a curated set of destination URL categories, collapsing all URLs pointing to the same Wikipedia page across languages into one. They do this classification in several ways, including incorporating the per-sample representation of the text span rather than merely the hyperlink token representations themselves (likely motivated by the likelihood of only a single hyperlink being present in the source text). We can realize this task as instances of several other common paradigms: (1) Single-task classification applied to the per-sample representation, (2) link prediction in a graph linking cross-lingual Wikipedia pages together, or (3) as an example of named entity recognition. This second stage is only allowed to modify the last two layers of the transformer architecture, which may be a vehicle to prevent catastrophic forgetting.
- [145] Sentiment Analysis; SAKG-BERT augments a pre-trained language model with a sentiment-analysis knowledge graph at the fine-tuning time only by concatenating relevant relationships from the KG based on sentiment-laden terms appearing in the review to the raw input text. They do not otherwise change the pre-training or fine-tuning process.