TABULAR LEARNING WITH BACKGROUND INFORMATION: LLMS, KNOWLEDGE GRAPHS, OR BOTH?

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026027028

029

031

033

034

037

040

041

042 043

044

046

047

048

051

052

ABSTRACT

Tables have their own structure, calling for dedicated tabular learning methods with the right inductive bias. These methods outperform language models. Yet, many tables contain text that refers to real-world entities, and most tabular learning methods ignore the external knowledge that such strings could unlock. Which knowledge-rich representations should tabular learning leverage? While large language models (LLMs) encode implicit factual knowledge, knowledge graphs (KGs) share the relational structure of tables and come with the promise of bettercontrolled knowledge. Studying tables in the wild, we assemble 105 tabular learning datasets comprising text. We find that knowledge-rich representations, from LLMs or KGs, boost prediction, and combined with simple linear models they markedly outperform strong tabular baselines. Larger LLMs provide greater gains, and refining language models on a KG boosts models slightly. On datasets where all entities are linked to a KG, LLMs and KG models of similar size perform similarly, suggesting that the benefit of LLMs over KGs is to solve the entity linking problem. Our results highlight that external knowledge is a powerful but underused ingredient for advancing tabular learning, with the most promising direction lying in the combination of LLMs and KGs.

1 Introduction: Background knowledge for Tabular Learning

Tabular learning Tabular data is central to machine-learning applications, powering applications from healthcare to finance. Yet, tables have properties that set them apart from other modalities. Cells may contain heterogeneous values: numbers, dates, categorical codes, or short texts. These values often only gain meaning through relational context, via column headers and neighboring entries. Tabular learning consists of making row-wise predictions, whether classification or regression, from these heterogeneous features. This unique structure has long favored learning methods with strong inductive biases for mixed-type features, such as gradient-boosted decision trees, over generic deep learning approaches (Grinsztajn et al., 2022; Shwartz-Ziv & Armon, 2022). Recent progress on table foundation models uses transformers with dedicated row-wise architectures, pretrained on synthetic tables (Hollmann et al., 2022; 2025), that are however purely numerical, leaving aside strings, dates, or categories. On the opposite, casting tables to text to readily apply large language models (LLMs) for in-context learning gives excellent few-shot performance, but does not scale nor benefit beyond a few dozen rows (Hegselmann et al., 2023; Gardner et al., 2024).

Text in tabular learners Tabular learners, unlike LLMs, leverage the specific repetitions of rows and features for state-of-the-art predictions on tables (Chen & Guestrin, 2016; Hollmann et al., 2025). Yet they also depart from LLMs in that they do not natively model text columns in tables. Here, a particularly underexplored dimension is that these texts often correspond to real-world entities, such as company names, drugs, or locations, that carry latent information far beyond the raw string. Exploiting this background knowledge could substantially improve prediction, especially in small-data regimes where tables themselves do not suffice to infer such knowledge from scratch. For example, a table of clinical trials mentioning drug names could benefit from external knowledge about drug classes, interactions, or approval status. However, state-of-the-art tabular learners are tailored to numbers (Erickson et al., 2025), using pipelines that cast entity strings to opaque numbers: categorical features are one-hot encoded, text reduced to surface-level representations, e.g., from character n-grams. Doing so discards the opportunity to ground table entries in external knowledge.

How can strings and entities bring background knowledge to tabular learning? A traditional answer would be to use data-integration and database techniques, augmenting tables with features obtained through joins with external databases (Doan et al., 2012; Cappuzzo et al., 2024). Yet, this approach faces well-known obstacles: discovering relevant tables, identifying joins, engineering relevant features while preventing their exponential growth across chained joins (Kanter & Veeramachaneni, 2015). A more scalable alternative enriches tables implicitly by mapping entity strings to vector representations pretrained on large-scale knowledge sources (Cvetkov-Iliev et al., 2023; Grinsztajn et al., 2023; Lefebvre & Varoquaux, 2025). Such embeddings provide compact summaries of factual and relational information from these sources, and can easily be injected into tabular models.

KGs and LLMs: two opposing philosophies of knowledge Pretraining embeddings from knowledge sources is shaped by two opposing philosophies of knowledge.

The *Knowledge Graph (KG) perspective* strives for pure, curated, knowledge. General-purpose KGs (Bollacker et al., 2008; Vrandečić & Krötzsch, 2014; Suchanek et al., 2024) gather facts in a structured, symbolic form, with a high signal-to-noise ratio: what they contain is largely correct. Their strength also lies in their explicit relational modeling, close to the relational nature of tabular data. Yet, their main weakness is their incompleteness: the number of true facts being potentially infinite, no KG can store them exhaustively.

By contrast, the *LLM perspective* treats knowledge as the statistical aggregation of written language. LLMs are probabilistic black-boxes trained on massive, weakly curated text corpora with no explicit notion of truth (Suchanek & Luu, 2023). They do not store curated facts, but model token co-occurrence statistics that implicitly encode fragments of factual knowledge (Petroni et al., 2019; Roberts et al., 2020; Jiang et al., 2020). Their power lies in breadth: scale enables coverage that far exceeds manually constructed KGs. This breadth comes at the price of reliability: LLMs are prone to hallucinations and factual drift (Ji et al., 2023; Tonmoy et al., 2024; Bang et al., 2025; Mallen et al., 2022), may produce confident but incorrect statements (Bender et al., 2021; Kadavath et al., 2022), and their internal reasoning remains opaque (Bender & Koller, 2020; Nanda et al., 2023). Raw application of LLMs to tabular learning also hits the wall of the size of their context window.

The knowledge integration bottleneck While LLMs can easily embed any string, the use of KGs in downstream tasks is hindered by a difficult knowledge integration step. Early KG embedding models operated in a transductive setting, learning representations for a fixed set of entities, primarily for internal tasks like link prediction (Bordes et al., 2013; Yang et al., 2014). Applying these embeddings to external data, such as tables, requires solving the challenging entity linking problem: mapping messy, real-world strings to canonical entities in the KG (Mendes et al., 2011; Delpeuch, 2019; Foppiano & Romary, 2020). This challenge is related to the broader "symbol grounding problem", a central difficulty of symbolic AI (Wikipedia, 2025). Recent advances in KG embedding models strive to overcome this limitation via generalization to unseen entities. One line of work couples KG embedding with pretrained (or jointly trained) text encoders applied to entity names or descriptions, so that unseen entities can be embedded directly from their textual mentions (Wang et al., 2021b; Saxena et al., 2022). A parallel effort focuses on building KG foundation models that can operate in a fully inductive setting, generalizing to entirely new graph structures by reasoning on their topology (Galkin et al., 2023; Huang et al., 2025a). These developments open up new avenues for integrating structured knowledge into downstream applications, but their effectiveness in the context of tabular learning remains an open question.

Contributions We study how to bring background information to tabular learning. Which modality, KGs or open-ended texts, should be preferred to pretrain world-knowledge models? Are numerical table foundation models all we need? What basic components for future research on table foundation models? To answer these questions, we assemble, from three diverse sources with different inclusion biases, 105 tabular learning datasets containing text. We conduct a large-scale empirical study, comparing, in a controlled setting, knowledge-rich representations from both LLMs and KG embedding models of varying sizes. We also study the impact of refining LLMs on KGs, to assess whether this hybrid approach combines the strengths of both modalities. Our findings are threefold:

1. **Bringing knowledge-rich representations into tabular learning matters:** both LLM and KG embeddings improve upon standard encoding techniques such as TF-IDF, bringing more gains than SOTA tabular learners developed for numerical tables.

- 2. **Refining LLMs on KGs is a promising combination:** clear gains are brought by scale, but refining LLMs on KGs also improves performance.
- 3. **Entity linking is the key bottleneck:** when all entities in a table are already linked to a KG, LLMs and KG models of comparable size perform similarly, suggesting that the main advantage of LLMs is their ability to implicitly solve the entity linking problem.

2 RELATED WORK

2.1 TABULAR LEARNING WITH TEXT FEATURES

From tree-based models to foundation models Historically, tabular learning has been dominated by gradient-boosted decision trees (GBDTs) (Chen & Guestrin, 2016; Ke et al., 2017; Prokhorenkova et al., 2018), which remain strong baselines due to their inductive biases for heterogeneous features (Grinsztajn et al., 2022). Recently, deep learning approaches (Ye et al., 2024; Holzmüller et al., 2024; Gorishniy et al., 2024), including table foundation models pretrained on synthetic data (Hollmann et al., 2022; 2025; Ma et al., 2024; Qu et al., 2025), now markedly outperform GBDTs (Erickson et al., 2025). However, a shared limitation of these methods is that they lack a dedicated mechanism for text features. Instead, they typically rely on simple string preprocessing turning these entries to numerical vectors, and then treat them as any other numerical feature. In practice, this vectorization step often ignores the semantics of string entries, relying on surface-level representations such as TF-IDF or character n-grams that bear no external knowledge.

Leveraging external knowledge from LLMs and KGs To address this gap, recent work has explored using external knowledge sources. One prominent approach is to leverage LLMs. Methods like TabLLM (Hegselmann et al., 2023) and Tabula-8B (Gardner et al., 2024) serialize table rows into text and fine-tune an LLM for classification and regression. These works put forward the benefit of in-context learning of LLMs, that brings their excellent few-shot performance to tabular learning, but cannot scale to the size of typical tables. Other work, such as TabStar (Arazi et al., 2025), adapt smaller, efficient text encoders with specialized architectures for tabular data. An alternative paradigm uses KGs as the source of external knowledge. For instance, CARTE (Kim et al., 2024) and TARTE (Kim et al., 2025) pretrain tabular models on KGs, but rely on the simple FastText (Bojanowski et al., 2017) model to process strings.

Prior comparative studies A few studies have begun to analyze the benefits of these knowledge-rich representations. Grinsztajn et al. (2023) demonstrated that embeddings from language models outperform traditional substring-based encoders, particularly for columns with diverse text entries. Similarly, Kasneci & Kasneci (2024) showed on 9 datasets that integrating embeddings from models like RoBERTa and GPT-2 into GBDTs often improves performance, especially in low-data regimes. While these works sketch out the value of using language models for text in tables, they do not inform of the relative merits of knowledge sourced from unstructured text (via LLMs) and structured graphs (via KG models).

2.2 Learning on KGs

Structure-based KG models A long-standing line of research learns representations from KGs by focusing solely on the graph structure. Early models operate in a transductive setting, learning low-dimensional embeddings for a fixed set of entities and relations. Such methods, that include TransE (Bordes et al., 2013), DistMult (Yang et al., 2014), ComplEx (Trouillon et al., 2016), and RotatE (Sun et al., 2019b), model the relations as geometric transformations in the embedding space and define a scoring function to measure the plausibility of triples. To overcome the limitations of transductive learning, subsequent work has focused on inductive models that can generalize to unseen entities (Zhu et al., 2021; Galkin et al., 2021). More recently, this has led to the development of KG foundation models that operate in a fully inductive setting, reasoning on the graph's topology to predict new links on entirely unseen graphs (Galkin et al., 2023; Lee et al., 2023; Huang et al., 2025a;b; Zhang et al., 2025; Du et al., 2025; Arun et al., 2025). Their application to tables remains however open, as it requires extracting from a table a relational graph rich-enough to enable the inductive setting.

Text-based KG models A parallel approach leverages the textual information associated with entities and relations, such as their names and descriptions. These models typically use a pretrained language model to create text-aware representations, bridging the gap between symbolic knowledge and natural language. One common strategy is to fine-tune a pretrained model such as BERT or RoBERTa using an objective that combines a masked language modeling loss with a KG-specific loss (Wang et al., 2021b;a; Youn & Tagkopoulos, 2022). Other methods frame link prediction as a textual task, either by scoring text sequences representing triples (Yao et al., 2019; Wang et al., 2022) or by treating it as a sequence-to-sequence problem where the model generates the missing entity's name (Chen et al., 2022; Xie et al., 2022). A prominent example of the latter is KGT5 (Saxena et al., 2022), which verbalizes triples and fine-tunes T5 (Raffel et al., 2020) to predict the missing elements. These text-based approaches enable embedding entities that were not seen during training, a crucial feature for downstream applications.

LLMs refined on knowledge Instead of training a model specifically for KG completion, another line of research refines general-purpose LLMs with structured knowledge to enhance their factual grounding. This approach aims to inject the high-quality, curated facts from KGs into the broader world knowledge implicitly stored in LLMs. For example, the ERNIE line of work (Sun et al., 2019a; 2020; 2021) refines language models like RoBERTa (Liu et al., 2019) by incorporating knowledge-base data into their pretraining objectives. More recently, the Knowledge Card framework (Feng et al., 2023) demonstrated that fine-tuning a moderately-sized LLM such as OPT-1.3B (Zhang et al., 2022) on KG triples can effectively plug factual knowledge into larger LLMs, improving their performance on knowledge-intensive tasks.

3 METHODOLOGY: A BENCHMARK FOR TABLE BACKGROUND KNOWLEDGE

3.1 105 TABULAR DATASETS

Three diverse data sources To ensure the robustness and generality of our findings, we assemble a diverse benchmark of 105 tabular datasets from three sources with distinct characteristics and inclusion biases: TextTabBench (Mráz et al., 2025), CARTE (Kim et al., 2024), and WikiDBs (Vogel et al., 2024).

TextTabBench and CARTE are established benchmarks for tabular learning, providing real-world tables with varied text features, from short entity names to longer descriptions. Each table is associated with a predefined prediction task (regression, binary, or multi-class classification). WikiDBs is a large corpus of over 1.6 million semi-synthetic tables generated from Wikidata. To create meaningful tasks from this source, we first filtered for tables with

Table 1: Task distribution across sources.					
Source	b-clf	m-clf	reg	Total	
TextTabBench	5	2	10	17	
CARTE	11	0	40	51	
WikiDBs	1	21	15	37	
Total	17	23	65	105	

Table 2: Aggregated features of tabular datasets across sources. The cardinality is computed on 1,024 rows.

TextTa	abBench	CARTE	WikiDBs
# columns	15.65	6.76	6.73
cardinality	286.36	371.44	463.70
string length	975.29	298.80	203.62
string similarity ¹	0.16	0.10	0.08

¹ cosine similarity of TF-IDF across rows

at least 1,200 rows, then manually curated a subset of 37 tables for which we could define a relevant prediction problem. Table 1 summarizes the final distribution of tasks across the three sources. Further details on each dataset are available in the Appendix (Table 6, Table 7, Table 8).

Data preprocessing We adopt the original preprocessing from TextTabBench and CARTE. For WikiDBs, we apply a procedure similar to TextTabBench. We also ensure that multi-class classification tasks have at most 10 classes, each with at least 105 samples. For all 105 datasets, we then apply the following preprocessing pipeline: (1) we remove all numerical columns to focus our study on text-based knowledge; (2) we log-transform regression targets with wide-ranging distributions; (3) we downsample majority classes in multi-class problems to create balanced datasets; and (4) we discard any table with fewer than 1,050 rows post-processing to ensure sufficient data for evaluation.

217

218

219

220

221

222

223

224

225

226

227

228229

230

231

232233234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249250

251

252

253 254

256

257

258

259260

261

262

263

264 265

266

267

268

Figure 1: An overview of our evaluation pipeline. For each dataset, we sample training and test sets. We then serialize the rows and use the embedding **model** to generate a vector representation for each row. Finally, we train a tabular learning estimator to evaluate these embeddings.

We also exclude one dataset from TextTabBench with excessively long text entries that exceed the context limits of some of our baselines.

Linked tables for controlled comparison

To isolate the contribution of knowledge from the challenge of entity linking, we create a specialized subset of 15 tables where text entries are unambiguously linked to entities in Wikidata5M (Wang et al., 2021b), a large-scale KG derived from Wikidata. These tables are selected from our main benchmark if they contain at least 1,050 rows with entities that can be matched to the KG. For this subset, we retain only the entity column and the prediction target, and re-

Table 3: Knowledge graph datasets. Smaller versions of Wikidata5M are created by filtering entities by degree ("deg."). All graphs use the same 822 relations.

	# entities	# triples	deg.
Wikidata5M	4.6M	20.6M	-
Wikidata3M	3.2M	15.5M	3
Wikidata2M	2.1M	11.5M	4
Wikidata1M	1.1 M	6.8M	6
Wikidata500k	0.5M	3.1M	9

move all unlinked rows. This setup allows for a direct comparison of pure KG models with LLMs in a scenario where entity linking is solved.

To analyze the impact of KG size, we generate four smaller KGs by progressively filtering out low-degree entities and retaining the largest connected component of the induced subgraph. The statistics of these graphs are presented in Table 3.

3.2 EVALUATION PIPELINE

Our evaluation pipeline, summarized in Figure 1, is designed to assess the quality of representations from various knowledge sources for downstream tabular tasks. For each dataset, we first generate embeddings from a given model and then train a tabular predictor on these embeddings to predict the target variable.

Experimental setup To simulate small-data scenarios where external knowledge is most critical, we sample training sets of varying sizes, $n_{train} \in \{64, 256, 1024\}$. The test set consists of 1,024 held-out samples (or all remaining samples if fewer are available). To ensure robust evaluation, we repeat this process 10 times with different random seeds for each configuration.

Embedding models We evaluate a wide range of models to generate representations, categorized as follows:

• Non-pretrained baseline: As a simple baseline without external knowledge, we use a TF-IDF vectorizer followed by a Truncated SVD with 30 components per column, implemented in the Skrub library.

- Pure LLMs: To study the effect of model scale and architecture, we include a diverse set of pretrained language models: the Llama-3.1 family (1B, 3B, 8B), the Qwen3 family (0.6B to 8B) (Yang et al., 2025) which performs well on the Massive Text Embedding Benchmark (Muennighoff et al., 2022), RoBERTa (base, large), T5 (small, base), and OPT-1.3B. We also include FastText as a representative of shallow, non-transformer text models.
- **Hybrid LLM+KG models:** To assess the benefit of structured knowledge, we evaluate models that refine LLMs on relational data. This includes ERNIE 2.0, KGT5, Knowledge Card, Tabula-8B, and TARTE. Each model is compared against its corresponding base LLM.
- Pure KG models: For the subset of 15 linked tables, we evaluate classic KG embedding models: DistMult, TransE, ComplEx, and RotatE. We train these models on Wikidata5M and its subsets, using an embedding dimension of d = 300.

Table serialization and downstream estimators To generate embeddings from LLMs, we serialize each table row into a natural language prompt. Following Gardner et al. (2024), we use the format: "The <col_a> is <val_a>. The <col_b> is <val_b>. What is the value of <target>?". For KGT5, we adapt the prompt to better match its pretraining format: "<col_a> | <val_a>. <col_b> | <val_b>. Predict: <target>". Contructing the embeddings across multiple columns (as opposed to the study of Grinsztajn et al. (2023)) is important because it enables the context (column name, other entries on the same row) to inform the representation, eg lead to disambiguate "Cambridge; UK" from "Cambridge; Massachusetts" in a table with columns "city; country".

The resulting high-dimensional embeddings are then fed into three representative tabular learners:

- Ridge regression: A simple and efficient linear model.
- **XGBoost:** A powerful GBDT model. To manage computational cost, we first reduce the embedding dimensionality to 300 using PCA. We then perform hyperparameter optimization via a randomized search (see Table 5).
- **TabPFNv2:** A transformer-based table foundation model, doing in-context learning. We use PCA to reduce dimensionality to 500, the maximum supported by the model.

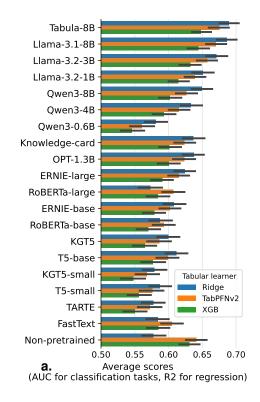
4 RESULTS: KNOWLEDGE REPRESENTATIONS FOR TABULAR LEARNING

4.1 Knowledge-rich representations boost tabular learning

More gains from knowledge representations than advanced tabular learning Figure 2 shows that improving the quality of the representations leads to more gains than using advanced tabular learning methods. Indeed, the best performance across the 105 datasets is obtained by a simple predictor, ridge, applied on good representations, such as those created via modern LLMs, outperforming sophisticated tabular learning methods XGBoost and TabPFNv2 (Figure 2a). In addition, more sophisticated tabular-learning models benefit less from advanced representations. This could be either because their flexibility enables them to fill-in for a less rich representation, or because the representations do not match their implicit inductive biases, tailored for tabular learning. Indeed, unlike typical tabular data, these representations are high-dimensional and closer to being rotationally-invariant Grinsztajn et al. (2022). Note that these advanced tabular learners cannot be applied as such to the knowledge-rich representations, as they have too many features for the corresponding implementations. Thus we need to reduce the input dimensions with PCA, following Grinsztajn et al. (2023).

A complementary observation is that the benefit of adding knowledge-rich representations to a simple tabular learner is larger than the benefit of using a sophisticated tabular learner on simpler representations: Figure 2b shows that TabPFNv2 achieves only half of the performance gains of Ridge combined with a good LLM-based representation.

Benefits for a wide variety of tables, from multiple sources Figure 2b shows that, for the ridge learner, knowledge-rich representations bring an improvement over non-pretrained string representation across methods, and larger models benefit consistently across the three different sources (Fig-



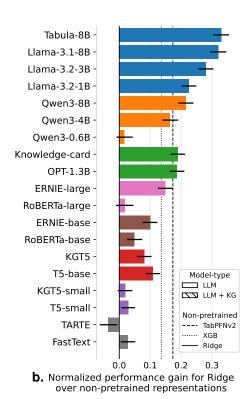


Figure 2: Performance gain of the various knowledge-rich representations compared to a non pretrained baseline – **a.** Comparisons including three tabular learners: ridge, XGBoost, and TabPFNv2; absolute scores. – **b.** Relative improvements to non-pretrained string representations, when using a ridge model as a tabular learner; normalized scores (0 is 10% worse, 1 is best score observed). – Appendix Figure 7 gives critical difference diagrams across all methods and datasets.

ure 8 gives source-specific results). These datasets are varied (Table 2), and the different sources represent different selections of tables with text. This diversity suggests that knowledge-rich representations help tabular learning in general, when the tables have text columns. The benefit is, on average, quite marked: going from non-pretrained string representation to the best LLM-based ones gives a .2 average boost in AUC or R2 to ridge (though only a .05 boost to TabPFNv2).

4.2 LARGER LLMS BRING MORE VALUE

Figure 3 shows the performance gain as a function of the LLM size (number of parameters), focusing only on pure LLM representations. It clearly reveals that the benefit increases as a function of size, for transformer-based representations (thus excluding FastText, which is a big model but very wide and shallow). This benefit of size is very clear in a given model family (comparing various sizes of e5, Qwen, or Llama-3). We hypothesize that this general scaling is driven by larger representational capacities brought by the increased number of parameters that enables the storage of more prior knowledge.

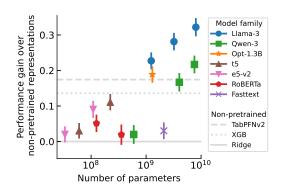


Figure 3: Effect of the size of the model, for pure-LLM representations.

4.3 REFINING LLMS ON KGS BOOSTS LANGUAGE MODELS SLIGHTLY

Figure 4 focuses on comparing the benefit brought by each method that has refined an LLM on a knowledge graph or knowledge base to the corresponding non-refined base LLM, as a function of size.

We estimate the scaling of the performance as a function of the number of parameters with a linear regression for both families of approaches –LLMs with and without KG refinement. Both families show the same scaling, but refining on KGs brings an offset: it enables reaching the same performance with a model with a number of parameters smaller by a factor of 2/3rd.

Note that the data points with the largest model correspond to the pair Tabula/Llama3 (Gardner et al., 2024), which refines on tabular data

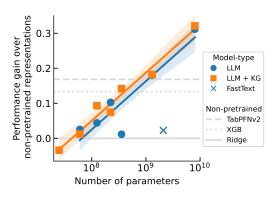


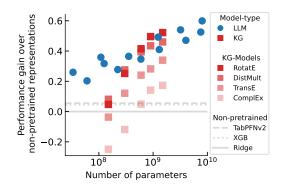
Figure 4: Comparison of LLM and their matched counterpart refined on knowledge bases.

rather than a rich KG. This pair also displays a comparatively smaller benefit of the refinement, which may result from the limited richness of the corresponding data.

The observed benefit of refining LLMs on KGs raises the question: what do knowledge graphs add to LLMs? How important is a rich knowledge graph?

4.4 TEASING OUT KNOWLEDGE FROM ENTITY MATCHING: TESTING PURE KG SOLUTIONS

LLMs are more than pure knowledge engineering objects: applied to embed texts, as we do here, they also bring in a form of fuzzy matching of entities (technically related to recontextualizing the tokens) and language understanding. This is to be contrasted with KGs, which are pure knowledge engineering objects (arguably with crisper knowledge), but 1) require entity matching and 2) do not bring language understanding. To tease out the role of background knowledge for tabular learning, we investigate a subset of tables for which the entity matching problem is solved, and each entry is linked to an entity in Wikidata5M.



In such an ideal scenario, pure KG embedding approaches provide features for the tables entries (Grover & Leskovec, 2016; Cvetkov-Iliev et al., 2023; Robinson et al., 2024). Figure 5

Figure 5: Comparing pure knowledge graphs to pure LLMs approaches on matched tables.

compares the benefits of LLM-based approaches with KG embedding approaches, varying the size of the models. For KG embedding, the size of the model is varied by varying the size of the KG used to build the embeddings (see Table 3): a smaller KG represents fewer entities, and thus has fewer parameters. When we reduce the size of the KG, it only provides representations for a fraction of the entities of the downstream table, and thus the downstream performance. This decrease is sharper than for LLMs, because smaller KGs face a hard failure (entity is not matched) while language models face a soft failure: they give an embedding whatever the query is. This embedding can be of varying quality, sometimes extrapolating beyond the knowledge of the LLMs, which corresponds to hallucination. However, an extrapolation that is only partly correct can still help downstream tabular learning.

Without entity-matching challenges, KG embedding is on par with LLMs For the largest, non-reduced KG, all table entries are matched, and good KG embedding models perform as well as LLMs

of the same size (Figure 5). Interestingly, this suggests that for the same number of parameters, KG embeddings do not store crispier knowledge than LLMs.

Driven by knowledge, rather than language understanding On the converse, when all entities are matched in the KG, LLMs of the same size do not bring a benefit, which suggests that for these tasks, the language-understanding features of LLMs are not important. Note that the selection of tables with entities all represented in KGs may introduce a selection bias towards more knowledge-centric tasks.

5 DISCUSSION AND CONCLUSION

Our large-scale study demonstrates that external knowledge is a powerful, yet underleveraged, ingredient for tabular learning. Representations from knowledge sources, whether LLMs or KGs, consistently improve prediction over standard text encodings. The gains from using better representations with a simple linear model surpass those from applying state-of-the-art tabular learners on less informative features, suggesting that for tables with text, the primary bottleneck is in representing this text well, rather than in the tabular learning algorithm.

LLMs solve symbol grounding, KGs provide curated knowledge Our results reveal a key trade-off. LLMs excel by implicitly solving the entity linking problem, mapping messy text to meaningful representations. Their performance scales with size, reflecting the vast knowledge encoded during pretraining. Conversely, pure KG models require explicit entity linking. Yet, when entities are prelinked, KG embeddings match LLMs of similar size. This implies, on the one hand, that the main advantage of LLMs here is not superior knowledge, but their ability to solve the symbol grounding problem by bridging unstructured text and canonical entities, and on the other hand, that pure-KG approaches, shying away from language models, do not bring the benefits of crispier knowledge for tabular learning.

A promising synergy: refining LLMs on KGs Our findings point to a potential synergy between LLMs and KGs. Refining LLMs on KGs improves performance, making models more parameter-efficient: a refined model achieves the performance of a pure LLM roughly 1.5 times its size.

Next-generation tabular foundation models should refine *large* language models. The clear benefit of scale points to a crucial direction for future research: developing large-scale, multimodal tabular foundation models pretrained on a combination of text and structured knowledge. Current tabular methods that jointly model strings and numbers use relatively small language models (Kim et al., 2024; 2025; Arazi et al., 2025), while our results stress the benefit of larger ones. The full potential of combining these large language models with KGs remains largely untapped. Massive, openly available knowledge bases like Wikidata, with more than 100M described entities, represent a rich yet underexploited resource for pretraining the next generation of tabular learners. In contrast, current methods pre-trained on knowledge (Kim et al., 2024; 2025) leverage only a small fraction of this available resource.

References

Mehdi Ali, Max Berrendorf, Charles Tapley Hoyt, Laurent Vermue, Sahand Sharifzadeh, Volker Tresp, and Jens Lehmann. Pykeen 1.0: a python library for training and evaluating knowledge graph embeddings. *Journal of Machine Learning Research*, 22(82):1–6, 2021.

Alan Arazi, Eilam Shapira, and Roi Reichart. Tabstar: A foundation tabular model with semantically target-aware representations. *arXiv preprint arXiv:2505.18125*, 2025.

Arvindh Arun, Sumit Kumar, Mojtaba Nayyeri, Bo Xiong, Ponnurangam Kumaraguru, Antonio Vergari, and Steffen Staab. Semma: A semantic aware knowledge graph foundation model. *arXiv* preprint arXiv:2505.20422, 2025.

- Yejin Bang, Ziwei Ji, Alan Schelten, Anthony Hartshorn, Tara Fowler, Cheng Zhang, Nicola Cancedda, and Pascale Fung. Hallulens: Llm hallucination benchmark. *arXiv preprint arXiv:2504.17550*, 2025.
 - Emily M Bender and Alexander Koller. Climbing towards nlu: On meaning, form, and understanding in the age of data. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pp. 5185–5198, 2020.
 - Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pp. 610–623, 2021.
 - Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the association for computational linguistics*, 5:135–146, 2017.
 - Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pp. 1247–1250, 2008.
 - Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems*, 26, 2013.
 - Riccardo Cappuzzo, Aimee Coelho, Felix Lefebvre, Paolo Papotti, and Gael Varoquaux. Retrieve, merge, predict: Augmenting tables with data lakes. *arXiv preprint arXiv:2402.06282*, 2024.
 - Chen Chen, Yufei Wang, Bing Li, and Kwok-Yan Lam. Knowledge is flat: A seq2seq generative framework for various knowledge graph completion. *arXiv preprint arXiv:2209.07299*, 2022.
 - Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794, 2016.
 - Alexis Cvetkov-Iliev, Alexandre Allauzen, and Gaël Varoquaux. Relational data embeddings for feature enrichment with background information. *Machine Learning*, 112(2):687–720, 2023.
 - Antonin Delpeuch. Opentapioca: Lightweight entity linking for wikidata. *arXiv preprint arXiv:1904.09131*, 2019.
 - AnHai Doan, Alon Halevy, and Zachary Ives. Principles of data integration. Elsevier, 2012.
 - Enjun Du, Siyi Liu, and Yongqi Zhang. Graphoracle: A foundation model for knowledge graph reasoning. *arXiv preprint arXiv:2505.11125*, 2025.
 - Nick Erickson, Lennart Purucker, Andrej Tschalzev, David Holzmüller, Prateek Mutalik Desai, David Salinas, and Frank Hutter. Tabarena: A living benchmark for machine learning on tabular data. *arXiv preprint arXiv:2506.16791*, 2025.
 - Shangbin Feng, Weijia Shi, Yuyang Bai, Vidhisha Balachandran, Tianxing He, and Yulia Tsvetkov. Knowledge card: Filling llms' knowledge gaps with plug-in specialized language models. *arXiv* preprint arXiv:2305.09955, 2023.
 - Luca Foppiano and Laurent Romary. entity-fishing: a dariah entity recognition and disambiguation service. *Journal of the Japanese Association for Digital Humanities*, 5(1):22–60, 2020.
 - Mikhail Galkin, Etienne Denis, Jiapeng Wu, and William L Hamilton. Nodepiece: Compositional and parameter-efficient representations of large knowledge graphs. *arXiv preprint* arXiv:2106.12144, 2021.
 - Mikhail Galkin, Xinyu Yuan, Hesham Mostafa, Jian Tang, and Zhaocheng Zhu. Towards foundation models for knowledge graph reasoning. *arXiv preprint arXiv:2310.04562*, 2023.

- Josh Gardner, Juan C Perdomo, and Ludwig Schmidt. Large scale transfer learning for tabular data
 via language modeling. Advances in Neural Information Processing Systems, 37:45155–45205,
 2024.
- Yury Gorishniy, Akim Kotelnikov, and Artem Babenko. Tabm: Advancing tabular deep learning with parameter-efficient ensembling. *arXiv* preprint arXiv:2410.24210, 2024.
 - Léo Grinsztajn, Edouard Oyallon, and Gaël Varoquaux. Why do tree-based models still outperform deep learning on typical tabular data? *Advances in neural information processing systems*, 35: 507–520, 2022.
 - Léo Grinsztajn, Edouard Oyallon, Myung Jun Kim, and Gaël Varoquaux. Vectorizing string entries for data processing on tables: when are larger language models better? *arXiv preprint arXiv:2312.09634*, 2023.
 - Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings* of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 855–864, 2016.
 - Stefan Hegselmann, Alejandro Buendia, Hunter Lang, Monica Agrawal, Xiaoyi Jiang, and David Sontag. Tabllm: Few-shot classification of tabular data with large language models. In *International conference on artificial intelligence and statistics*, pp. 5549–5581. PMLR, 2023.
 - Noah Hollmann, Samuel Müller, Katharina Eggensperger, and Frank Hutter. Tabpfn: A transformer that solves small tabular classification problems in a second. *arXiv preprint arXiv:2207.01848*, 2022.
 - Noah Hollmann, Samuel Müller, Lennart Purucker, Arjun Krishnakumar, Max Körfer, Shi Bin Hoo, Robin Tibor Schirrmeister, and Frank Hutter. Accurate predictions on small data with a tabular foundation model. *Nature*, 637(8045):319–326, 2025.
 - David Holzmüller, Léo Grinsztajn, and Ingo Steinwart. Better by default: Strong pre-tuned mlps and boosted trees on tabular data. *Advances in Neural Information Processing Systems*, 37:26577–26658, 2024.
 - Xingyue Huang, Pablo Barceló, Michael M Bronstein, Ismail Ilkan Ceylan, Mikhail Galkin, Juan L Reutter, and Miguel Romero Orth. How expressive are knowledge graph foundation models? *arXiv preprint arXiv:2502.13339*, 2025a.
 - Xingyue Huang, Mikhail Galkin, Michael M Bronstein, and İsmail İlkan Ceylan. Hyper: A foundation model for inductive link prediction with knowledge hypergraphs. *arXiv preprint arXiv:2506.12362*, 2025b.
 - Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM computing surveys*, 55(12):1–38, 2023.
 - Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438, 2020.
 - Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
 - James Max Kanter and Kalyan Veeramachaneni. Deep feature synthesis: Towards automating data science endeavors. In 2015 IEEE international conference on data science and advanced analytics (DSAA), pp. 1–10. IEEE, 2015.
 - Gjergji Kasneci and Enkelejda Kasneci. Enriching tabular data with contextual llm embeddings: A comprehensive ablation study for ensemble classifiers. *arXiv preprint arXiv:2411.01645*, 2024.
 - Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, 30, 2017.

- Myung Jun Kim, Leo Grinsztajn, and Gael Varoquaux. Carte: Pretraining and transfer for tabular learning. In *ICML*, 2024.
- Myung Jun Kim, Félix Lefebvre, Gaëtan Brison, Alexandre Perez-Lebel, and Gaël Varoquaux. Table foundation models: on knowledge pre-training for tabular learning. *arXiv preprint arXiv:2505.14415*, 2025.
- Jaejun Lee, Chanyoung Chung, and Joyce Jiyoung Whang. Ingram: Inductive knowledge graph embedding via relation graphs. In *International Conference on Machine Learning*, pp. 18796–18809. PMLR, 2023.
- Félix Lefebvre and Gaël Varoquaux. Scalable feature learning on huge knowledge graphs for downstream machine learning. *arXiv* preprint arXiv:2507.00965, 2025.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Junwei Ma, Valentin Thomas, Rasa Hosseinzadeh, Hamidreza Kamkari, Alex Labach, Jesse C Cresswell, Keyvan Golestan, Guangwei Yu, Maksims Volkovs, and Anthony L Caterini. Tabdpt: Scaling tabular foundation models. arXiv preprint arXiv:2410.18164, 2024.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. *arXiv* preprint arXiv:2212.10511, 2022.
- Pablo N Mendes, Max Jakob, Andrés García-Silva, and Christian Bizer. Dbpedia spotlight: shedding light on the web of documents. In *Proceedings of the 7th international conference on semantic systems*, pp. 1–8, 2011.
- Martin Mráz, Breenda Das, Anshul Gupta, Lennart Purucker, and Frank Hutter. Towards benchmarking foundation models for tabular data with text. *arXiv preprint arXiv:2507.07829*, 2025.
- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*, 2022.
- Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023.
- Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. Language models as knowledge bases? *arXiv preprint arXiv:1909.01066*, 2019.
- Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. Catboost: unbiased boosting with categorical features. Advances in neural information processing systems, 31, 2018.
- Jingang Qu, David Holzmüller, Gaël Varoquaux, and Marine Le Morvan. Tabicl: A tabular foundation model for in-context learning on large data. *arXiv preprint arXiv:2502.05564*, 2025.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the parameters of a language model? *arXiv preprint arXiv:2002.08910*, 2020.
- Joshua Robinson, Rishabh Ranjan, Weihua Hu, Kexin Huang, Jiaqi Han, Alejandro Dobles, Matthias Fey, Jan Eric Lenssen, Yiwen Yuan, Zecheng Zhang, et al. Relbench: A benchmark for deep learning on relational databases. *Advances in Neural Information Processing Systems*, 37:21330–21341, 2024.
- Apoorv Saxena, Adrian Kochsiek, and Rainer Gemulla. Sequence-to-sequence knowledge graph completion and question answering. *arXiv preprint arXiv:2203.10321*, 2022.

- Ravid Shwartz-Ziv and Amitai Armon. Tabular data: Deep learning is not all you need. *Information Fusion*, 81:84–90, 2022.
 - Fabian Suchanek and Anh Tuan Luu. Knowledge bases and language models: Complementing forces. In *International Joint Conference on Rules and Reasoning*, pp. 3–15. Springer, 2023.
 - Fabian M Suchanek, Mehwish Alam, Thomas Bonald, Lihu Chen, Pierre-Henri Paris, and Jules Soria. Yago 4.5: A large and clean knowledge base with a rich taxonomy. In *Proceedings of the 47th international ACM SIGIR conference on research and development in information retrieval*, pp. 131–140, 2024.
 - Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, and Hua Wu. Ernie: Enhanced representation through knowledge integration. *arXiv* preprint arXiv:1904.09223, 2019a.
 - Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. Ernie 2.0: A continual pre-training framework for language understanding. In *Proceedings of the AAAI* conference on artificial intelligence, volume 34, pp. 8968–8975, 2020.
 - Yu Sun, Shuohuan Wang, Shikun Feng, Siyu Ding, Chao Pang, Junyuan Shang, Jiaxiang Liu, Xuyi Chen, Yanbin Zhao, Yuxiang Lu, et al. Ernie 3.0: Large-scale knowledge enhanced pre-training for language understanding and generation. *arXiv preprint arXiv:2107.02137*, 2021.
 - Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. Rotate: Knowledge graph embedding by relational rotation in complex space. *arXiv preprint arXiv:1902.10197*, 2019b.
 - SMTI Tonmoy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. A comprehensive survey of hallucination mitigation techniques in large language models. *arXiv* preprint arXiv:2401.01313, 6, 2024.
 - Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. Complex embeddings for simple link prediction. In *International conference on machine learning*, pp. 2071–2080. PMLR, 2016.
 - Liane Vogel, Jan-Micha Bodensohn, and Carsten Binnig. Wikidbs: A large-scale corpus of relational databases from wikidata. Advances in Neural Information Processing Systems, 37:41186–41201, 2024.
 - Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85, 2014.
 - Bo Wang, Tao Shen, Guodong Long, Tianyi Zhou, Ying Wang, and Yi Chang. Structure-augmented text representation learning for efficient knowledge graph completion. In *Proceedings of the Web Conference 2021*, pp. 1737–1748, 2021a.
 - Liang Wang, Wei Zhao, Zhuoyu Wei, and Jingming Liu. Simkgc: Simple contrastive knowledge graph completion with pre-trained language models. *arXiv* preprint arXiv:2203.02167, 2022.
 - Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhengyan Zhang, Zhiyuan Liu, Juanzi Li, and Jian Tang. Kepler: A unified model for knowledge embedding and pre-trained language representation. *Transactions of the Association for Computational Linguistics*, 9:176–194, 2021b.
 - Wikipedia. Symbol grounding problem Wikipedia, the free encyclopedia. http://en.wikipedia.org/w/index.php?title=Symbol%20grounding%20problem&oldid=1291979022, 2025. [Online; accessed 25-September-2025].
 - Xin Xie, Ningyu Zhang, Zhoubo Li, Shumin Deng, Hui Chen, Feiyu Xiong, Mosha Chen, and Huajun Chen. From discrimination to generation: Knowledge graph completion with generative transformer. In *Companion proceedings of the web conference* 2022, pp. 162–165, 2022.
 - An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

language models. arXiv preprint arXiv:2205.01068, 2022.

- Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases. arXiv preprint arXiv:1412.6575, 2014. Liang Yao, Chengsheng Mao, and Yuan Luo. Kg-bert: Bert for knowledge graph completion. arXiv preprint arXiv:1909.03193, 2019. Han-Jia Ye, Huai-Hong Yin, and De-Chuan Zhan. Modern neighborhood components analysis: A deep tabular baseline two decades later. arXiv e-prints, pp. arXiv-2407, 2024. Jason Youn and Ilias Tagkopoulos. Kglm: Integrating knowledge graph structure in language models for link prediction. arXiv preprint arXiv:2211.02744, 2022. Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christo-
 - Yucheng Zhang, Beatrice Bevilacqua, Mikhail Galkin, and Bruno Ribeiro. Trix: A more expressive model for zero-shot domain transfer in knowledge graphs. *arXiv preprint arXiv:2502.19512*, 2025.

pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer

Zhaocheng Zhu, Zuobai Zhang, Louis-Pascal Xhonneux, and Jian Tang. Neural bellman-ford networks: A general graph neural network framework for link prediction. *Advances in neural information processing systems*, 34:29476–29490, 2021.

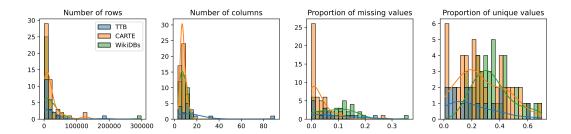


Figure 6: Statistics distribution across sources.

A MORE DETAILS ON THE EXPERIMENTS

More statistics on datasets Figure 6 gives statistics about table sizes, proportion of missing values, and mean column cardinality.

Table 6, Table 7 and Table 8 provide details on each individual dataset.

Table 4: Task distribution across sources, for linked tables.

Source	b-clf	m-clf	reg	Total
CARTE	0	0	3	3
WikiDBs	1	3	8	12
Total	1	3	11	15

Experiments on linked tables We have 15 linked tables, 4 for classification and 11 for regression. Details on these tables are provided in Table 4.

For the KG embedding models (DistMult, TransE, ComplEx and RotatE), we use d=300 for the embedding dimension, and train them for 100 epochs with a batch size of 8192 and a learning rate of 10^{-3} , and use the default parameters of their PyKEEN implementation (Ali et al., 2021).

Metrics and score normalization We evaluate performance using the R2 score for regression and the ROC-AUC score for classification. To aggregate results across datasets of varying difficulty, we normalize scores for each dataset and random seed. Following Grinsztajn et al. (2022), we establish a normalized scale where the best-performing model scores 1 and the model at the 10th performance percentile scores 0. Other models' scores are mapped to this [0, 1] range via an affine transformation. For regression, we clip scores at 0 to mitigate the impact of poor-performing outliers.

Uncertainty estimation To account for statistical variability, we repeat each experiment 10 times with different random seeds. The error bars in our result figures represent the standard error to the mean across these runs.

XGBoost hyperparameter tuning For the XGBoost estimator, we perform hyperparameter optimization via a randomized search with 100 iterations. We use 5-fold cross-validation, repeated 5 times on the training set, to evaluate each hyperparameter configuration. The detailed search space is provided in Table 5.

B ADDITIONAL RESULTS

B.1 RUNTIME ANALYSIS

The benefits of leveraging external knowledge come at a computational cost. Table 9 details the average runtimes for embedding generation and estimator fitting (ridge) across different embedding models and training sizes. As expected, larger models introduce a significant computational overhead. For instance, generating embeddings with an 8-billion-parameter LLM is, on average, over 100 times slower than using the non-pretrained baseline. This highlights the trade-off between predictive performance and the computational resources required for knowledge integration.

Table 5: Search space for XGBoost hyperparameters.

Hyperparameter	Distribution	Range
n_estimators	Integer	[50, 1000]
max_depth	Integer	[2, 6]
min_child_weight	Log-uniform	[1, 100]
subsample	Uniform	[0.5, 1.0]
learning_rate	Log-uniform	$[10^{-5}, 1]$
colsample_bylevel	Uniform	[0.5, 1.0]
colsample_bytree	Uniform	[0.5, 1.0]
gamma	Log-uniform	$[10^{-8}, 7]$
reg_lambda	Log-uniform	[1, 4]
alpha	Log-uniform	$[10^{-8}, 100]$

822 823 824

825

826 827

Table 6: Overview of TextTabBench datasets used in our benchmark. Table statistics after preprocessing.

columns

5

12

10

13

16

9

4

33

6

14

14

9

17

5

13

90

13

classes

2

2

2

2

2

4

10

rows

17,000

1,732

18,720

11.254

3,598

1,384

10,000

3,818

2,914

11,349

1,165

1.253

12,000

183,794

19,427

1,281

37,484

Task

b-clf

b-clf

b-clf

b-clf

b-clf

m-clf

m-clf

reg

linked rows

837

843

848 849 850

851

852 853

854 855

856

857

858

859

860

861

862

863

B.2 OVERALL MODEL RANKING

Dataset

Diabetes

Spotify

Airbnb

Beer

Job Frauds

Kickstarter

Lending Club

Osha Accidents

Customer Complaints

California Houses

Insurance Complaints

San Francisco Permits

Stack Overflow

Covid Trials

IT Salary

Mercari

Wine

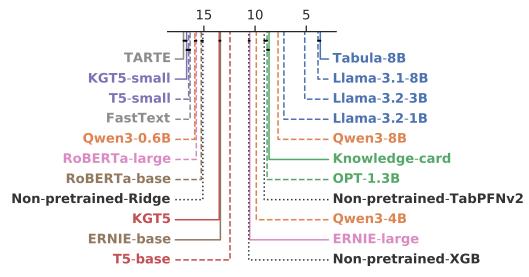
Figure 7 presents a critical difference diagram comparing the mean ranks of all embedding methods when paired with a Ridge predictor. It also includes the performance of more advanced estimators on non-pretrained representations for context.

B.3 PERFORMANCE ANALYSIS BY DATA SOURCE

Figure 8 illustrates the relative improvements of knowledge-rich representations over non-pretrained ones, broken down by data source. The benefits of external knowledge vary with dataset characteristics; tables from WikiDBs and CARTE, which are more knowledge-intensive, gain more from these representations than those from TextTabBench.

Figure 9 details the effect of LLM size on performance for each data source, confirming the scaling trend across different types of tables.

Figure 10 compares the performance of base LLMs to their counterparts refined on KGs. The benefits of refinement are most pronounced for the WikiDBs datasets, which are inherently more knowledge-centric as they are derived from a knowledge base.



Average rank measured at n = 1024

Figure 7: Critical difference diagram across all data sources and methods.

C LLM USAGE

LLMs were used to polish the writing of some parts of this paper.

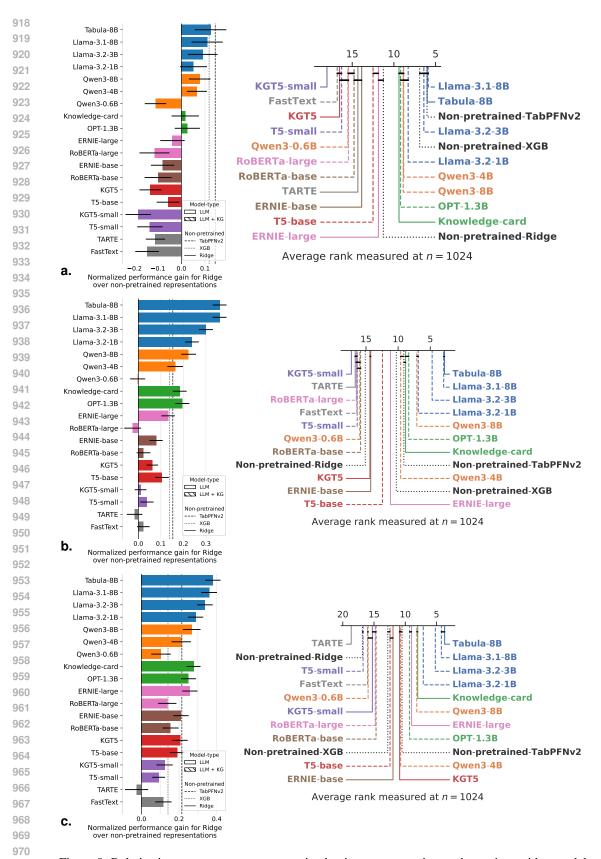


Figure 8: Relative improvements to non-pretrained string representations, when using a ridge model as a tabular learner for each source Larger models consistently yield better performances.: **a.** Text-TabBench **b.** CARTE **c.** WikiDBs.

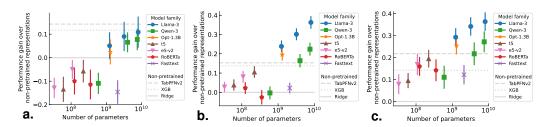


Figure 9: Effect of the size of representations from pure-LLM models for each source: **b.** CARTE **c.** WikiDBs.

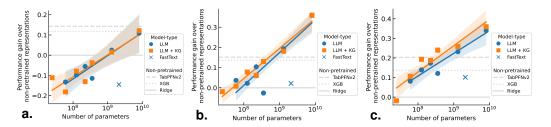


Figure 10: Comparison of LLM and the matched models refined on knowledge bases for each source: **b.** CARTE **c.** WikiDBs.

Table 7: Overview of CARTE datasets used in our benchmark. Table statistics after preprocessing.

Dataset	Task	# rows	# columns	# classes	# linked rows
Chocolate Bar Ratings	b-clf	2,218	7	2	-
Coffee Ratings	b-clf	1,670	9	2	-
Michelin	b-clf	6,774	6	2	-
NBA Draft	b-clf	1,550	5	2	-
Ramen Ratings	b-clf	3,726	5	2	-
Roger Ebert	b-clf	2,668	6	2	-
Spotify	b-clf	41,096	8	2	-
US Accidents Severity	b-clf	20,930	10	2	-
Whisky	b-clf	1,788	7	2	-
Yelp	b-clf	60,088	9	2	-
Zomato	b-clf	60,302	8	2	-
Movies	reg	7,224	8	-	7,095
US Accidents Counts	reg	22,623	7	-	14,697
US Presidential	reg	19,857	7	-	13,221
Anime Planet	reg	14,391	7	-	, -
Babies R Us	reg	5,085	5	_	_
Beer Ratings	reg	3,197	6	_	_
Bikedekho	reg	4,786	6	_	_
Bikewale	reg	8,992	6	_	_
Buy Buy Baby	reg	10,718	5	_	<u>-</u>
Cardekho	reg	37,813	14	_	_
Clear Corpus	reg	4,724	11	_	_
Company Employees	reg	10,941	8		
Employee Remuneration	reg	35,396	3	_	-
Employee Salaries	_	9,211	7	_	-
Fifa22 Players	reg		10	-	-
-	reg	18,085		-	-
Filmty Movies	reg	41,205	7 5	-	-
Journal JCR	reg	9,615		-	-
Journal SJR	reg	27,931	10	-	-
Japanese anime	reg	15,535	12	-	-
K-Drama	reg	1,239	9	=	=
ML/DS salaries	reg	10,456	8	-	=
Museums	reg	11,467	15	-	=
Mydramalist	reg	3,400	11	-	-
Prescription Drugs	reg	1,714	6	-	-
Rotten Tomatoes	reg	7,158	11	-	-
Used Cars 24	reg	5,918	7	-	=
Used Cars Benz Italy	reg	16,391	6	-	=
UsedCars.com	reg	4,009	9	-	-
Used Cars Pakistan	reg	72,655	5	-	-
Used Cars Saudi Arabia	reg	5,507	8	-	-
Videogame Sales	reg	16,410	5	-	-
Wikiliq Beer	reg	13,461	8	-	=
Wikiliq Spirit	reg	12,275	6	-	-
Wina Poland	reg	2,247	13	_	-
Wine.com Prices	reg	15,254	7	-	-
Wine.com Ratings	reg	4,095	7	_	-
WineEnthusiasts Prices	reg	120,975	9	_	-
WineEnthusiasts Ratings	reg	129,971	9	_	_
WineVivino Price	reg	13,834	6	_	_

Table 8: Overview of WikiDBs datasets used in our benchmark. Table statistics after preprocessing.

Dataset	Task	# rows	# columns	# classes	# linked rows
CC Authors	b-clf	16,224	8	2	1,302
Defenders	m-clf	18,610	11	10	8,700
Philosophers	m-clf	4,230	9	10	1,656
US Music Albums	m-clf	3,270	11	10	2,180
Artist Copyrights	m-clf	2,000	10	10	-
Artworks Catalog	m-clf	1,210	9	10	-
Forward Players	m-clf	1,400	11	10	-
Geographers	m-clf	1,130	10	10	-
Historic Buildings	m-clf	27,980	7	10	-
Islands	m-clf	19,650	4	10	-
Kindergarten Locations	m-clf	2,790	7	3	-
Magic Narratives	m-clf	1,062	5	9	-
Museums	m-clf	9,550	5	10	-
Noble Individuals	m-clf	1,400	10	10	-
Notable Trees	m-clf	1,408	5	8	-
Parish Churches	m-clf	1,350	5	10	-
Sculptures	m-clf	3,720	7	10	-
Spring Locations	m-clf	5,930	3	10	-
State Schools	m-clf	2,800	4	10	-
Scientific Articles	m-clf	2,760	14	10	-
Sub Post Offices	m-clf	1,530	4	10	-
Transport Stations	m-clf	4,640	9	10	-
Business Locations	reg	16,821	5	-	16,438
Dissolved Municipalities	reg	13,462	7	-	1,656
Geopolitical Regions	reg	1,114	7	-	1,066
Historical Figures	reg	11,260	12	-	2,134
Municipal District Capitals	reg	1,658	6	-	1,267
Poets	reg	60,240	11	-	21,564
Territorial Entities	reg	36,717	8	-	34,189
WWI Personnel	reg	30,675	12	-	16,227
Artworks Inventory	reg	10,635	6	-	-
Drawings Catalog	reg	63,130	9	-	-
Eclipsing Binary Stars	reg	297,934	7	-	-
Registered Ships	reg	4,644	7	-	-
Research Articles	reg	6,962	7	-	-
Research Article Citations	reg	4,115	10	-	-
Ukrainian Villages	reg	21,355	4	-	-

Table 9: Average runtimes for embedding extraction and ridge fitting.

	Train-size			
Models	64	256	1 024	
Tabula-8B Llama-3.1-8B Llama-3.2-3B Llama-3.2-1B	123.86 ± 141.32 119.44 ± 133.49 43.22 ± 50.78 18.07 ± 20.18	145.15 ± 165.82 140.00 ± 156.66 50.77 ± 59.65 21.34 ± 23.69	215.92 ± 257.74 209.29 ± 246.60 75.76 ± 92.71 31.69 ± 36.55	
Qwen3-8B Qwen3-4B Qwen3-0.6B	119.54 ± 143.78 64.99 ± 82.19 12.01 ± 12.62	140.11 ± 168.76 76.23 ± 96.49 14.23 ± 14.75	209.66 ± 262.17 114.36 ± 150.38 20.99 ± 21.96	
Knowledge-card	25.48 ± 29.22	30.16 ± 34.27	44.85 ± 53.11	
OPT-1.3B	22.82 ± 28.97	27.06 ± 33.98	40.37 ± 52.91	
ERNIE-large	8.28 ± 5.87	10.00 ± 6.77	14.50 ± 9.00	
RoBERTa-large	7.92 ± 5.59	9.60 ± 6.51	14.09 ± 9.31	
ERNIE-base	4.79 ± 3.76	5.94 ± 4.30	8.35 ± 5.51	
RoBERTa-base	4.22 ± 3.05	5.28 ± 3.47	7.43 ± 4.26	
KGT5	4.60 ± 3.43	5.80 ± 3.98	8.32 ± 5.33	
T5-base	5.31 ± 6.26	6.55 ± 7.29	9.24 ± 10.69	
KGT5-small	3.22 ± 2.91	4.22 ± 3.31	5.60 ± 3.84	
T5-small	3.50 ± 3.00	4.42 ± 3.43	5.94 ± 3.96	
TARTE	4.44 ± 4.33	5.37 ± 5.01	7.75 ± 6.28	
FastText	2.30 ± 3.80	3.05 ± 4.42	3.85 ± 5.64	
Non-pretrained-Ridge	0.52 ± 0.65	1.27 ± 1.21	2.14 ± 2.26	