

Technical Report - SAILS: Schema Naming Assessments for Improved LLM-Based SQL Inference

Kyle Luoma

kluoma@ucsd.edu

University of California, San Diego

La Jolla, California, USA

Arun Kumar

akk018@ucsd.edu

University of California, San Diego

La Jolla, California, USA

ABSTRACT

Large Language Models (LLMs) have revolutionized Natural Language to SQL (NL-to-SQL), dominating most NL-to-SQL benchmarks. But LLMs still face limitations due to hallucinations, semantic ambiguity, and lexical mismatches between an NL query and the database schema. Naturally, a lot of work in the ML+DB intersection aims to mitigate such LLM limitations. In this work, we shine the light on a complementary data-centric question: *How should DB schemas evolve in this era of LLMs to boost NL-to-SQL?* The intuition is that more NL-friendly schema identifiers can help LLMs work better with DBs. We dive deeper into this seemingly obvious, but hitherto underexplored and important, connection between schema identifier “naturalness” and the behavior of LLM-based NL-to-SQL by creating a new integrated benchmark suite we call SAILS. SAILS has 4 novel artifacts: (1) A collection of real-world DB schemas not present in prior NL-to-SQL benchmarks; (2) A set of labeled NL-SQL query pairs on our collection not seen before by public LLMs; (3) A notion of naturalness level for schema identifiers and a novel labeled dataset of modified identifiers; and (4) AI artifacts to automatically modify identifier naturalness. Using SAILS, we perform a comprehensive empirical evaluation of the impact of schema naturalness on LLM-based NL-to-SQL accuracy, and present a method for improving LLM-based NL-to-SQL with natural views. Our results reveal statistically significant correlations across multiple public LLMs from OpenAI, Meta, and Google on multiple databases using both zero-shot prompting as well as more complex NL-to-SQL workflows: DIN SQL, and CodeS. We present several fine-grained insights and discuss pathways for DB practitioners to better exploit LLMs for NL-to-SQL.

CCS CONCEPTS

- Information systems → Data management systems; Database design and models; Database performance evaluation; Data cleaning.

KEYWORDS

Database, Relational database schema, schema design, natural language to SQL, LLM, benchmark

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGMOD/PODS '25, June 03–05, 2025, Berlin, Germany

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

ACM Reference Format:

Kyle Luoma and Arun Kumar. 2025. Technical Report - SAILS: Schema Naming Assessments for Improved LLM-Based SQL Inference. In *Proceedings of The 2025 ACM SIGMOD/PODS Conference (SIGMOD/PODS '25)*. ACM, New York, NY, USA, 41 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Natural language-to-SQL (NL-to-SQL) query generation capability has been revolutionized by foundational large language models (LLMs) [5, 36, 48]. This has made the integration of LLM-based query tools into relational database workflows more viable, with both established DBMS vendors and startups beginning to offer commercial NL-to-SQL interfaces. However, challenges in the NL-to-SQL space remain that can degrade the effectiveness of an LLM-enabled data retrieval workflow in real-world databases [15]. Principal among such challenges is *schema linking*, which is the association of entities in NL utterances to elements in the database schema.

While much work has studied making LLMs larger or more sophisticated, a more basic issue often underlies this challenge: lexical mismatches between natural language and poorly-named tables and columns in a schema. Intuitively, schema elements that are “better named” could raise the accuracy of schema linking within the NL-to-SQL setup. In this paper, we unpack and dive deeper into this intuition to study how exactly the “naturalness” of schema elements matters for NL-to-SQL by instituting a new benchmark and performing extensive empirical analysis using that. One might ask: *Why bother formalizing a concept that seems obvious and intuitive?* We believe this is important for 2 reasons. First, without a more formalized—or at least automated way—to define, verify, and compare “naturalness” researchers and practitioners alike will be forced to grapple with ad hoc and inconsistent approaches. In turn, this can lead to confounded conclusions by researchers on how different LLMs behave on different schemas and mislead practitioners comparing different NLIs. This points to the need for a new benchmark labeled dataset for this problem.

Second, practitioners need a way to efficiently and accurately operationalize any insights about the impact of naturalness on their schema elements for LLM-based NLIs. This points to the need for a systematic evaluation of how naturalness affects different databases, queries, and LLMs used for NL-to-SQL.

Our Focus. In this paper, we take the first steps toward deeper understanding on this seemingly obvious, but hitherto underexplored and important, relationship between schema identifier naturalness and LLM-based NL-to-SQL. Specifically, we ask the following three interconnected questions. (1) How do we quantify “naturalness” of schema identifiers? (2) Does it really impact schema

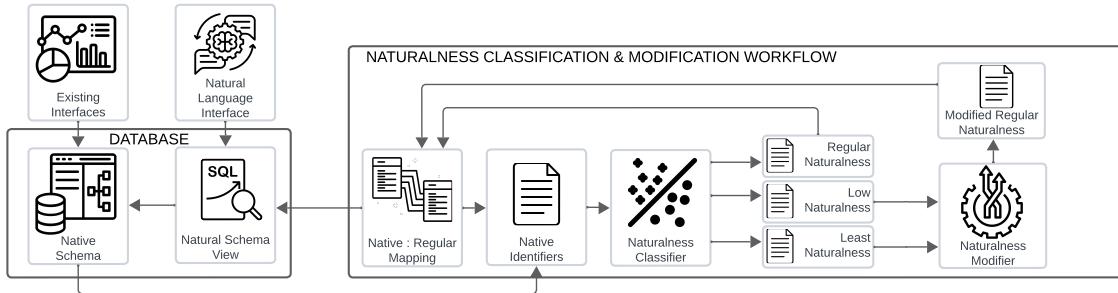


Figure 1: Databases with poorly named, or less natural, schema identifiers perform poorly in LLM-based NL-to-SQL interfaces, and this project exposes the need for more natural schemas. We offer approaches and artifacts, including a naturalness classification and modification workflow, that can aid in the naturalness assessment and modification processes required to create a performance-enhancing natural view. In this way, the native schema remains as-is so that existing tools can continue talking to it without modification, while an LLM-based NLI can be integrated into the existing stack via a natural view.

linking accuracy in LLM-based NL-to-SQL and if so, by how much? (3) How does that impact vary by complexity of the database and query, as well as across different popular LLMs?

To answer the above questions, we create a novel integrated benchmark suite we call SAILS with new collections of real-world databases and query pairs, a new labeled dataset of schema identifiers, a set of evaluation metrics, and LLM prompting and other AI artifacts.

1.1 Preliminaries and Setup

LLM-based NL-to-SQL. The most obvious way to seek LLM performance improvements would be by increasing the power of the language models themselves. But the cost of training and deploying LLMs continues to increase in concert with their complexities. Additionally, many practitioners seek “plug and play” solutions by employing already-available LLMs. Model training and finetuning impose access barriers that may render such a pursuit untenable for organizations that use databases but lack the requisite talent such as data science and machine learning expertise.

The practice of prompt engineering can also help improve NL-to-SQL performance, though dealing with schema complexity and schema representations in LLM prompting is an ongoing challenge in enterprise-level NL-to-SQL applications [15]. The majority of leading submissions on the popular Spider NL-to-SQL benchmark leaderboard are LLM-based solutions [12, 17, 46] that employ a variety of prompting strategies, some of which require multiple successive API requests containing schema context and instructions. These approaches can be costly and unintuitive for NLDB end users, and can incur excessive costs and overhead when deployed at scale.

A complementary line of work on realistic NL-to-SQL benchmarking uses structural schema modification such as normalization, flattening, and replacement to evaluate effects on LLM performance. Making such structural changes to target schemas challenges model robustness and increases error rates in NL-to-SQL performance [31], and this recent work indicates that schema design is a viable target for LLM-based NL-to-SQL accuracy improvements.

Schema Linking. Schema linking remains as a persistent challenge for LLMs. With the availability of capable LLMs that consistently generate valid SQL statements, a larger proportion of

NL-to-SQL generation errors are now associated with incorrect or ambiguous database identifier selection as opposed to incorrect syntax [52]. Schema linking performance has been improved using lexical matching heuristics [20, 62], joint relationally aware embeddings with attention [6, 58], the use of pre-trained language models to perform schema probing [59], and multimodel pipelines with ML models for pruning schema knowledge [27]. Some NL-to-SQL methods address schema linking challenges by adding additional context such as sample values or metadata [46] to schema knowledge representations. These methods can improve performance in some cases [34], and can be useful for schemas with obscurely-named tables and columns, though they do so at the cost of much larger schema knowledge representations.

Schema linking still often fails, even with the most capable LLMs due to poorly-aligned schema identifier names with natural language question contents, that could be due to the use of synonyms or the obscurity of a database identifier. In the latter case, it can be challenging for even a sophisticated linking solution to match natural language words to schema elements that yield minimal semantic meaning.

Schema Naming Conventions. The majority of database schema naming best practices originate from *practitioners* and are generally published as software documentation, organization policies, tutorials, etc. We find that there is a gap in database and data integration *academic* literature evaluating schema identifier naming practices for any purpose. While the semantics of schema identifiers may not have been considered as a necessary subject of database research in the past, the increasing integration of natural language interfaces to databases has elevated its importance.

Naming conventions for database schema identifiers vary by organization, database vendor, application, and purpose. A web search for database table and column naming guidelines yields multiple resources ranging from blog posts [7], StackOverflow responses [50], DBMS vendor documentation [39], and tutorials [19]. Poor schema identifier naming practices is considered a database code smell [49] where meaningless identifier names should be avoided. Generally, the most consistent best practices include selecting descriptive and

concise names that contain only commonly-understood abbreviations and acronyms, though some conventions suggest the use of abbreviated prefix and suffix modifiers that describe application associations, or entity purpose [40].

In our research, we identified several databases containing schemas with varying levels of human-readability and understandability (what we will call naturalness) which suggests that there can be a tendency for database schema designers to choose shorter and less descriptive identifier naming conventions. As we will see, such naming shortcuts can negatively affect NL-to-SQL performance.

1.2 Our Benchmark Artifacts and Analyses

Given the above context of our benchmarking setup, we now explain the new artifacts in SAILS, followed by a summary of our empirical analysis.

Artifact 1: Real-World Database Schemas. The SAILS benchmark contains several new *real-world database schemas* that are not part of existing NL-to-SQL benchmarks (Artifact 1). Our focus on schema naming motivates the creation of a new novel benchmark dataset, because existing benchmark naturalness levels are higher than those of many real-world schemas, and other real-world schema collections including SchemaPile [11] lack the necessary database instances to enable NL-to-SQL evaluation. In our analysis of these real-world schemas, we discover that identifier naming variances generally appear in the form of abbreviations and expansions; we refer to these variances as identifier *naturalness*.

Artifact 2: Identifier Naturalness Classifications. Our analysis reveals that naturalness can be formalized categorically with the help of finetuned language models and feature engineering. We then hand-label the schema identifiers, with some ML assistance, to classify their naturalness level and produce a new golden labeled dataset. We classify identifiers into one of 3 naturalness levels (Regular, Low, and Least) (Artifact 2). This dataset, consisting of over 17,000 labeled identifiers, serves as the training data for the naturalness classifiers described next.

Artifact 3: Naturalness Classifiers. We experiment with various classification approaches, and make available the models trained to classify the naturalness of a database schema identifier (Artifact 3).

Artifact 4: Naturalness-Modified Identifiers. To better understand the effect of schema identifier naturalness, and to enable within-database experiments, we create alternate versions of each real-world schema identifier at each naturalness level (Artifact 4). This dataset serves two purposes: 1) Training data for ML-based naturalness modifiers, and 2) Generation of schemas with varying naturalness levels to analyze the impact of naturalness on NL-to-SQL performance. We modify the identifiers with the assistance of LLM prompting, finetuned models, and database metadata.

Artifact 5: Naturalness Modifier. We offer an in-context learning-based prompting strategy for identifier naturalness reduction (or abbreviation). We also provide an identifier naturalness increaser (or expander) that leverages retrieval augmented generation, interactive few-shot example building, and database metadata parsing methods to streamline the database naturalness improvement process.

Artifact 6: NL-to-SQL Question Query Pairs. The SAILS benchmark contains 503 NL question-SQL query pairs which we use for NL-to-SQL performance analysis of 4 LLMs. We created this new collection as another hand-labeled golden dataset without the use of AI-based workflows (Artifact 6).

Experimental Evaluation. Using the SAILS benchmark artifacts, we analyze and experiment with the effects of schema identifier naturalness on LLM NL-to-SQL performance. We select 5 publicly-available LLMs: OpenAI’s GPT-3.5, GPT-4o, a finetuned variant of Meta’s Code-Llama, Google’s newest Gemini 1.5, and CodeS finetuned for NL-to-SQL. We evaluate them using both execution result set matching and a novel identifier set comparison approach that pinpoints schema linking performance.

In this paper we focus primarily on a simple zero-shot prompting of the LLM for our experiments. We recognize that this may not be the best for overall execution accuracy, but it helps us isolate the impact of schema identifier naturalness in this first work on this problem. As such, more complex workflows will create confounding effects while not necessarily providing more insights into schema linking performance. However, for completeness sake, we also compare two illustrative complex workflows: DIN SQL for task-specific prompt chaining [46], and CodeS [28] for NL-to-SQL finetuning.

We find that schema identifier naturalness by and large does have a meaningful effect on NL-to-SQL accuracy and schema linking performance. Specifically, identifier naturalness is moderately and positively correlated with both schema linking and execution accuracy. Identifiers of low naturalness yield lower performing NL-to-SQL inferences in terms of both schema linking (identifier recall) and execution accuracy. These findings have implications for practitioners who are either designing new databases intended for LLM-based applications, or seeking to augment existing RDBMSs with an LLM-based NL-to-SQL interface.

In summary, this paper makes the following contributions:

- We propose a novel measure of *naturalness* of a database schema identifier and demonstrate through extensive experiments that naturalness has a significant effect on LLM schema linking performance in the context of NL-to-SQL.
- We provide a hybrid LLM-generated and human-curated training dataset (Artifact 2) and language model (Artifact 3) for schema naturalness classification.
- We offer a new multi-domain NL-to-SQL evaluation benchmark collection consisting of 9 real-world relational databases (Artifact 1) and 503 unpublished NL-to-SQL query pairs (Artifact 6) that do not exist in any LLM training corpora.
- We create a novel labeled dataset of alternate naturalness levels that map the identifiers from Artifact 1 to hybrid LLM-human curated identifiers of different naturalness levels (Artifact 4), and methods for expanding and abbreviating identifiers to change their naturalness (Artifact 5).
- We conduct an extensive empirical analysis of the performance of 5 popular foundational LLMs over our benchmark using a novel schema linking metric for NL-to-SQL.
- We propose a realistic workflow that enables the preservation of existing database integrations while offering LLM-based NLIs a natural view of a target schema.

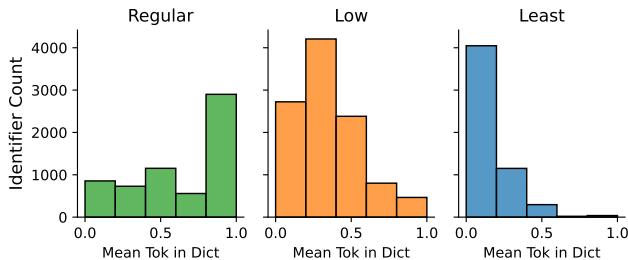


Figure 2: Mean Token in Dictionary, the proportion of tokens in an identifier that match a word in an English dictionary, generally aligns with the SAILS 3-class naturalness categorization approach.

Regular	Low	Least
airbag	AccountChk	AdCtTxIRWT
AdaptiveCruiseControl	IsueFrDate	COGM_Act
ModelYear	RecvAsst	DfltlSp
service_name	UsrQuery	FNDAbs
Research_Staff	ValueOfT	CSI22

Table 1: Example identifiers and their naturalness levels, from the SAILS naturalness labeled dataset (Artifact 2).

2 SCHEMA IDENTIFIER NATURALNESS

Intuitively, naturalness can be thought of as the degree to which a phrase, or word, resembles natural language. Naturalness is a concept and target of research in field of controlled natural languages [25], where controlled language syntax is evaluated in terms of naturalness levels. Recent NL-to-SQL research also defines and measures naturalness [31] for the purpose of evaluating the naturalness of natural language question utterances, but avoids measuring the naturalness of schema elements.

To the best of our knowledge, no prior attempts have been made to definitively measure the naturalness of a database schema's identifiers. In order to achieve this goal, we propose a three-category naturalness classification scheme in order to measure the effects of naturalness on NL-to-SQL performance.

2.1 Naturalness Categories

As the first work on this topic of how schema identifier naturalness affects LLMs, we seek to define a preliminary metric—one that is consistent and descriptive enough to differentiate between naturalness levels and to measure their effects.

To gain insights into naturalness-related trends in the SAILS datasets, we create a *mean token-in-dictionary* measurement that describes the proportion of tokens in an identifier that exactly match a word in a comprehensive English word list. Figure 2 reveals differences between each naturalness category where Least naturalness identifiers contain fewer in-dictionary tokens, and Regular naturalness identifiers are more likely to consist of in-dictionary tokens. This distribution suggests that because the bulk of the training corpora of LLMs is human-generated natural language text, what humans consider “natural” for such identifiers generally aligns with how LLMs react to them.

Examples of schema identifiers and their naturalness categories are displayed in Table 1. We define these categories with the underlying assumption that the identifiers are named as some semantic representation of the data, and that naming-related problems of interest are related how an identifier is codified. That is, identifiers are assumed to not be random character sequences or random words that do not correspond to the content of the database entities they represent. With this assumption in mind, we categorize naturalness into 3 discrete levels as follows:

- Regular: The identifier contains complete English words with no abbreviations or acronyms, or contains only acronyms in common usage (e.g., ID or GPS).
- Low: The identifier contains abbreviated English words and less common acronyms that are usually recognizable by non-domain experts (e.g., UTM or CPI). The meaning of the identifier can be inferred without consulting external documentation.
- Least: The identifier’s meaning cannot be inferred by non-experts due to indecipherable acronyms or abbreviations, and external metadata or other documentation must be consulted in order to determine its purpose.

While we recognize that naturalness can also be treated as a continuous spectrum, between the choices of continuous scoring and discrete categories, we select the latter as an initial approach to naturalness evaluation. The primary factors underlying this choice are the level of effort required to conduct human-based scoring of a large set of database identifiers, and the difficulty of consistently scoring naturalness on a continuous range over a large set of data. Therefore, we use an intuitive and easily-verifiable discrete 3-class taxonomy in the first work on this topic.

2.2 Naturalness Classification

To consider naturalness as a factor in NL-to-SQL performance, we derive naturalness scores of the target schemas’ identifiers. We use this score to consider effects of individual identifier naturalness, schema naturalness, and query identifier naturalness. Because manual naturalness classification can be a time consuming task for large schemas, we automate the process by training a machine learning-based classifier. This effort is beneficial in multiple situations. First, it can ease some manual effort of the labeling process and make the process of scaling to more databases in the future less labor intensive. Second, it can help practitioners efficiently and consistently evaluate the naturalness of their own database schema identifiers prior to NLI integration.

To train a classifier to perform identifier naturalness scoring, we employ the 3-class set of naturalness categories described in Section 2.1, and a list of database identifiers drawn from the SAILS real-world database schemas (Artifact 1). We categorize the naturalness of each identifier to generate the SAILS *identifier naturalness classification* labeled data (Artifact 2) which we use for ML-based naturalness classifier training, evaluation and testing.

We evaluate multiple classification approaches including heuristic-based word matching, few-shot LLM prompting with GPT-3.5 and GPT-4, and LLM finetuning. The GPT-4 few-shot approach achieves 74 percent accuracy and an f1 score of .77. We experiment with multiple finetuning collections, first using a hand-labeled collection of

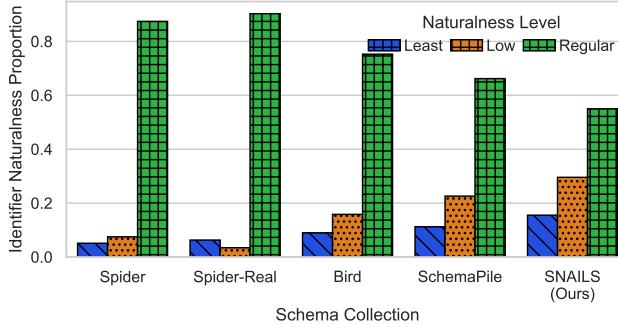


Figure 3: Comparison of the SNAILS database collection (Artifact 1) described in Section 3.1 to other real-world and benchmark schema collections. SNAILS naturalness proportions are generally biased toward less natural identifiers and is more consistent with the real-world SchemaPile collection than other existing benchmarks including Spider and Spider Realistic.

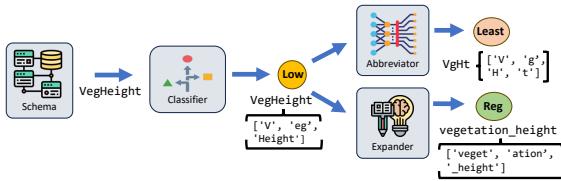


Figure 4: Schema identifiers are classified (Artifact 2) and modified to increase or decrease naturalness as appropriate. Modified identifiers comprise the schema crosswalks used for schema modification during NL-to-SQL experimentation (Artifact 3).

1,648 naturalness classifications and then leveraging the initial classifier along with weak supervision to generate a larger collection of 17,226 labeled identifiers. Finetuning using the second collection outperforms all few-shot approaches, with the two best-performing classifiers fine-tuned GPT 3.5 and BERT-based CANINE [8] models performing at 89 percent accuracy, and 0.89 f1 score.

Figure 3 provides a visual comparison between the SNAILS schema collection and common NL-to-SQL benchmarks including Spider, Spider Realistic, and BIRD. Additionally, we compare the SNAILS collection to the real-world SchemaPile collection and find that SNAILS collection proportions generally align to SchemaPile naturalness, more so than other existing benchmarks, which creates a more realistic and challenging benchmark in terms of schema naturalness.

To better understand the magnitude of naming practices in real-world schemas, we use the CANINE-based classifier to classify the naturalness of the SchemaPile collection: a large volume of real-world database schemas [11] that contains over 22,000 database schemas, 198,000 tables, and 1 million columns. We find that in over 7,500 schemas (32 percent of the collection) Least natural identifiers make up at least 10 percent of the schema identifier names.

Additionally, over 5,000 schemas register a combined naturalness of 0.7 or below—an indicator that the schema contains a high level of Low and Least naturalness identifiers. We examined the naturalness category distribution for these 5,000 schemas, and found that for this subset of schemas Low and Least naturalness identifiers outnumber Regular naturalness identifiers. These findings reinforce the importance of the naturalness problem by confirming that, although a reasonable majority of schemas are already natural, there still exist many schemas with lower naturalness levels in the real-world—enough to motivate the formalization of schema naming quality measures.

2.3 Identifier Schema Naturalness Mapping

In addition to measuring the effects of identifier naturalness in existing schemas, we also seek to evaluate the effects of modifying schema naturalness. For this purpose, we create Artifact 4, naturalness-modified identifiers. This artifact enables schema modification during prompt generation and query inference, which provides a within-schema assessment of naturalness level effects on NL-to-SQL accuracy.

Identifier Mapping. In addition to the ground truth, or Native, naturalness of the 9 schemas in the SNAILS real-world database collection, the naturalness-modified identifier collection contains 3 additional sets of identifiers: Regular, Low, and Least. That is, each native identifier is mapped to 2 additional, semantically equivalent, identifiers of higher or lower naturalness, and mapped to itself for its own naturalness level (i.e., we do not generate new identifiers of the same naturalness as its native form).

Figure 4 provides a visual example of the Native identifier *VegHeight* which is classified as Low naturalness. With this naturalness classification as a starting point, we abbreviate it further to generate a corresponding Least naturalness identifier *VgHt*. We also expand it to generate the corresponding Regular naturalness version *vegetation_height*. We map the Native *VegHeight* identifier to itself in the Low naturalness category.

Naturalness Modification. For *more natural to less natural* modifications (the abbreviator in Figure 4), we employ in-context learning (few-shot) prompt strategies with GPT-3.5 turbo to generate naturalness-modified identifiers (e.g., Regular to Low, Low to Least, and Regular to Least). We favor this approach over model finetuning, as simple instructions to abbreviate the identifier coupled with several examples prove more effective and less prone to poor results (e.g., presence of unwanted characters in the modified identifier).

Automating the reverse *less natural to more natural* naturalness modification (the expander in Figure 4) requires additional context and external knowledge from data description sources. Though a recent project describes a promising identifier expansion strategy [63] without external knowledge, it requires finetuning over a large dataset, and is likely susceptible to overfitting; therefore we opt for our own approach that incorporates the use of an LLM augmented with schema metadata lookup capability. To accomplish this, we create a Python program with GPT interaction that takes as input metadata describing a schema's native identifiers, and outputs

Database	Tables	Columns	Questions	Org
ASIS	36	245	40	NPS
ATBI	28	192	40	NPS
CWO	13	71	40	NPS
KIS	18	157	40	NPS
NPFM	27	190	40	NPS
NTSB	40	1611	100	NHTSA
NYSED	27	423	63	NYSED
PILB	21	196	40	NPS
SBOD	2588	90,477	100	SAP

Table 2: SAILS Real-World Database Schemas

an identifier with regular naturalness. More details of this process are available in the appendix.

3 BASE COLLECTIONS

Given the recency of the LLMs selected for evaluation in this project, and the relative maturity of existing NL-to-SQL benchmarks, we believe that foundational LLMs have been exposed to existing benchmark training and development NL questions and queries in their training corpora. NL-to-SQL performance differences between queries over seen vs. unseen schema are significant [52], and we seek to avoid as much bias as possible due to intentional or unintentional pre-training on existing benchmark datasets.

We also find that existing benchmarks including Spider [61], and BIRD [29], do not match the identifier naturalness distribution of real-world schema collections such as SchemaPile [11]. Although SchemaPile is a very large representation of real-world schemas, it does not contain database instances necessary for benchmark performance evaluations; so, we are not able to leverage its dataset in the creation of a new benchmark. To reduce bias due to benchmark data exposure, and to create a benchmark more representative of real-world schema naming, SAILS contains two artifacts for NL-to-SQL benchmarking: Artifact 1, which is a collection of 9 publicly-available database schemas and data; and Artifact 6, a human-generated set of 503 NL question - gold query pairs.

3.1 Datasets

Native Schemas. The SAILS real-world database schema collection (Artifact 1) consists of 9 databases sourced from multiple locations. We refer to the schema identifier names as they exist in the source databases as *Native*, and we classify each schemas' Native naturalness level (see Figure 5). Domain diversity facilitates a more thorough evaluation [14]; so, SAILS database collections span multiple domains. Domain coverage includes scientific nature observation records, vehicle safety statistics, primary school performance data, and business resource planning.

The U.S. National Parks Service's IRMA Portal [1] is the source of the scientific observation databases which include the Field Data for the Inventory of Amphibians and Reptiles of Assateague Island National Seashore (**ASIS**) [9], Great Smoky Mountains All Taxa Biodiversity Inventory (**ATBI**) Plot Vegetation Monitoring Database [13], Wildlife Observations Database: Craters of the Moon National Monument and Preserve 1921-2021 (**CWO**) [51], Exotic and Invasive Plants Monitoring Database (**KIS**) [23], Northern Plains

Fire Management (**NPFM**) [33] and Pacific Island Network Landbird Monitoring Dataset (**PILB**) [24].

The National Transportation Safety Bureaus 2021 safety sampling dataset [35, 47] is the source of SAILS **NTSB** safety statistics database. We source school performance data (**NYSED**) from the New York State Education Department [4].

The business resource planning database **SBOD** is a training example of the popular SAP Business One system, and is publically available in MS SQL server backup format [45]. The SBOD schema consists of an extremely large number of tables and columns; so pruning is required to fit the schema within the context window of the LLMs we compared. We reduce the schema knowledge token requirements by segmenting the SBOD schema into submodules and further reducing tables through data profiling. Additional information on the SBOD schema knowledge management is available in the appendix.

Each database was migrated from its source format into an MS SQL Server database. Several databases contained identifiers with whitespace characters, which is uncommon in most schemas. To mitigate whitespace-related inference failures as a confounder, we modify the native identifiers by replacing whitespace characters with underscore characters. In total, 148 out of over 19,000 total identifiers (less than .01 percent) contained at least 1 whitespace character.

Native Schema Naturalness Levels. Since the intent of this project is to measure the effect of schema naturalness, we first check if there is sufficient distribution of naturalness levels across the collection. We employ the GPT-3.5-based classifier described in Section 2.2 to evaluate the naturalness of the native schema identifiers.

In addition to measuring the proportion of identifiers in each naturalness category, we also derive a combined naturalness score. Combined naturalness is the weighted average of category proportion values, where scores range from 0.0 to 1.0 with 1.0 representing a schema containing only Regular naturalness identifiers. A more detailed description of its calculation is available in the appendix.

Figure 5 displays the proportions of identifiers in each naturalness category, as well as the combined naturalness, in each native schema. From the chart, we can see that the schemas in the SAILS collection described in Section 3.1 represent a heterogeneous selection of naturalness combinations.

Modified (Virtual) Schemas. To control for confounding factors such as schema structure, normalization levels, and constraint variances between native schemas, we perform within-database evaluations of naturalness. To accomplish this, we generate 3 additional *virtual* schemas using the naturalness-modified identifiers (Artifact 4) described in Section 2.3. Each virtual schema is representative of a naturalness category, where schema identifiers are replaced with a semantically equivalent identifier of a different naturalness level. This results in 4 schema versions per database in the base collection: Native, Regular, Low, and Least.

The modified schemas are virtual because we do not create database instances that can be queried directly. Rather, we query virtual schemas via identifier replacement in prompts and generated queries using processes described in Section 4. This approach reduces storage overhead. It also enables possible future schema

Database	Qs	Top	Function	Join	CK Join	Exists	Subquery	Where	Negation	Group By	Order By	Having
ASIS	40	1	24	13	1	0	2	18	0	17	1	0
ATBI	40	5	20	18	0	1	7	21	2	16	7	1
CWO	40	2	18	5	1	5	10	34	7	12	2	1
KIS	40	8	26	15	0	0	2	25	1	11	8	0
NPFM	40	5	27	21	0	0	1	29	0	16	5	0
NTSB	100	8	82	23	21	0	6	62	4	42	23	4
NYSED	63	10	36	10	4	1	21	55	1	16	10	1
PILB	40	6	25	23	0	0	3	20	0	16	11	2
SBOD	100	2	33	44	0	0	0	82	0	17	2	1

Table 3: Gold query clause counts for each SAILS database. Columns represent a count of gold queries that contain the listed clause types. Qs is the count of question-query pairs for each database. CK Join is the subset of joins that require a composite key. Note: MS SQL Server dialect replaces the common LIMIT clause with an equivalent TOP clause that precedes select items in the SELECT clause.

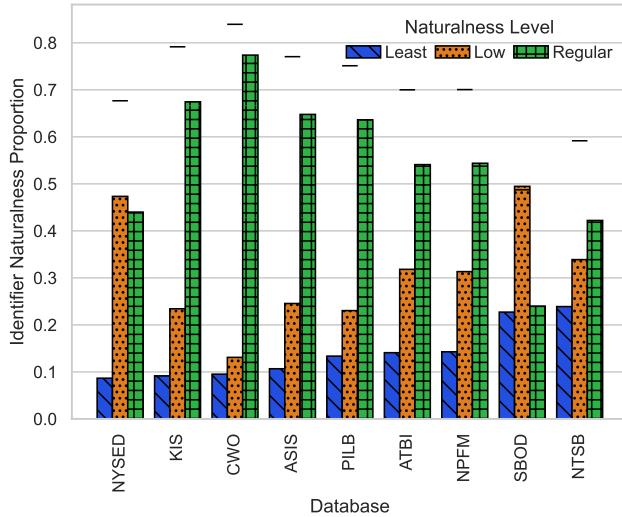


Figure 5: Proportion of identifiers in each naturalness category within the SAILS real-world database collection (Artifact 1). Horizontal line markers indicate calculated combined naturalness as described in the appendix

variations of different naturalness proportions without the need to instantiate additional database instances.

SAILS Database Selection and Extension Processes. The initial 9 datasets and schemas are included because they were (1) publicly available, (2) not included in any prior NL-to-SQL benchmarks, (3) contained relational tables with dependencies and database instances with values, (4) had available table and column metadata, (5) represented a diversity of application domains, and (6) contain data potentially useful for real-world data analysis or data science applications. Databases are not selected or pre-screened using perceived naturalness as criteria.

We view the initial 9 schemas as a starting point from which the SAILS dataset can grow. Researchers who wish to extend the SAILS collection should use the same selection criteria. In addition,

the extension process must ensure that new databases: (1) can be represented as MS SQL Server instances, (2) each native identifier's naturalness is classified according to defined criteria using the SAILS naturalness classifier, and (3) that native identifiers are modified using the SAILS modification artifacts to create alternate naturalness levels.

3.2 NL Question - SQL Query Pairs

To evaluate SQL inference performance over the Native and modified schemas in the SAILS real-world database collection, we create a new set of 503 NL-question and SQL gold query pairs (Artifact 6). Schema identifier naturalness are the primary considerations for NL question and gold query composition. During question and query formulation we track schema coverage to ensure that the distribution of identifier naturalness within a set of gold queries generally matches the naturalness distribution of target schemas.

To enable accuracy measurements at the identifier level, gold queries contain the minimum identifiers (tables and columns) required to answer its corresponding question. For this reason, for questions that require the count aggregation function, where appropriate, we use the COUNT(*) clause (as opposed to selecting an arbitrary column). This approach eliminates incorrect penalties to recall if a generated query fails to project an arbitrary column as a function argument.

Gold queries contain only native identifiers, such that all gold queries return valid non-null results from target databases in the real-world database collection (Artifact 1). We measure query complexity as a count of its clauses and identifiers. Gold queries span a range of complexities, from very simple single table projections, to multi-table joins and nested subqueries (see Table 3).

Adding New NL-SQL Pairs to the SAILS Collection. For researchers interested in extending the SAILS collection, it is necessary to create new ground truth NL-SQL pairs for evaluation. While we employed a fully manual approach for question writing, and this approach may be used for future additions, they may also consider the use of new approaches such as using a template-based approach for generating question-query pairs with relational data as input [42]. Regardless of NL-SQL pair creation method, researchers

should ensure adequate schema coverage and minimum essential identifier selection as described in the preceding section.

4 NL-TO-SQL BENCHMARKING SETUP

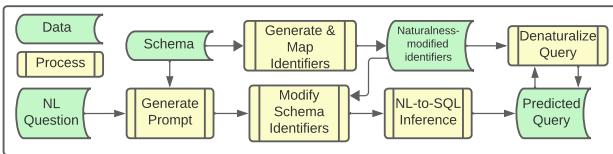


Figure 6: Experiment setup workflow from NL question and schema as input to predicted query as output.

To evaluate the impact of naturalness on NL-to-SQL accuracy, we build a benchmarking setup pipeline as shown in Figure 6. NL question and gold query pairs, database schemas, and naturalness crosswalk mappings are inputs into subprocesses. The subprocesses include prompt generation, schema identifier naturalness modification, identifier naturalness classification, LLM-based NL-to-SQL inference, and predicted query “denaturalization” (i.e., converting table and column identifiers to native schema identifiers prior to query execution). The output of the experiment setup is a predicted query, which along with its gold query counterpart, is executed against a target database. This predicted query is passed into a parser analysis tool as initial steps of the *Performance Evaluation and Results Classification* phase of the experiment described in Section 5.

4.1 Prompt Generation

The design space for LLM-based NL-to-SQL prompting is quite large, with options ranging from zero-shot instructions to sequential prompting broken into discrete tasks such as schema subsetting and error handling. Although we evaluate 2 complex NL-to-SQL workflows, to maintain consistency across the LLMs compared in this study, our performance comparisons focus on a single prompting strategy: zero-shot prompting with schema knowledge.

Prompting Strategy. SAILS prompts consist of zero-shot instructions with schema knowledge (denoted as ZS in results figures) in a format similar to OpenAI’s Text-to-SQL demonstration prompt [17] for completions. The prompt begins with task instructions and database information:

For the database described next, provide only a sql query.
do not include any text that is not valid SQL.

#Database: NTSB

#MS SQL Server tables, with their properties:

Target database system tables provide schema knowledge, which is represented as a list of tables and their column names with data types in the format:

#TableName (Col1Name Type, Col2Name Type, ...)

The prompt ends with the instruction:

a sql query, written in the MS SQL Server dialect,
to answer the question: <Question>

Where <Question> is replaced with an NL question directed at the given schema.

To evaluate naturalness effects on more complex NL-to-SQL prompting workflows, we also implement DIN SQL [46] which uses prompt chaining with GPT-4, and CodeS [28]—a multi-step NL-to-SQL system (schema filtering and SQL inference) based on StarCoder [30] and finetuned for the NL-to-SQL translation task.

Prompt Schema Identifier Modification. For inference on virtual schemas with modified naturalness levels, we replace Native identifiers with corresponding identifiers of the target virtual schema’s naturalness level. We accomplish this step using the naturalness-modified identifier collection (Artifact 4) described in Section 2.3. We use a SQL parser to encase identifiers within tags to improve identifier replacement accuracy and eliminate errors due to substring matching between identifiers.

4.2 NL-to-SQL Inference

Language Models. Foundational LLMs continue to grow in capability at a rapid pace. Despite this growth, not all NLI implementations can avail of the most-capable LLMs, often due to organizational policy constraints (e.g., organizational security concerns [18]). Additionally, we seek to understand if schema naming effects generalize across model architectures and sizes. Thus, we consider several LLMs, both open and closed source, to capture as many use profiles as possible including OpenAI’s GPT-3.5 Turbo and GPT-4o [36, 37]; Google’s Gemini 1.5 Ultra [53, 54]; and Phind-CodeLlama-34B-v2 [44] which is a finetuned variant of Meta’s CodeLlama 2 [48].

CodeS and DIN SQL Implementation. For the more complex DIN SQL and CodeS NL-to-SQL workflows, we provide additional versions of the SAILS schema artifacts to conform to the input requirements of the target systems. Additionally, we add data logging between agents to document the schema filtering step for additional analysis. For consistency between approaches, we use GPT-4o for all steps in the prompting chain. For CodeS inference, we execute the schema filtering and NL-to-SQL inference using the CodeS codebase and finetuned models.

Generated Query Denaturalization. For queries targeted at virtual schemas and generated using modified schema identifiers, we perform reverse modifications prior to query execution on the native database schema. Using a purpose-built Antlr [43]-based parser, we extract table and column identifiers, and generate a tagged query with identifier tags encasing table and column names. The tags guide the replacement algorithm, ensuring accurate replacement of naturalness-modified identifiers with their Native naturalness counterparts.

5 NL-TO-SQL BENCHMARKING RESULTS

This section describes the process of evaluating the generated SQL query output from the prior section. We evaluate performance in terms of execution accuracy (result set comparison and manual evaluation) and schema linking (recall, precision, and F1).

Key Takeaways. Overall, there is a model-dependent statistically significant correlation between identifier naturalness and

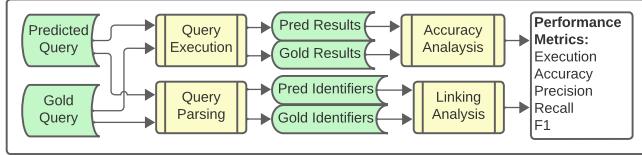


Figure 7: Benchmark results evaluation includes generated and gold query execution on target schemas, parser-based analysis, and identifier set comparisons. We evaluate performance in terms of execution accuracy and schema linking (precision, recall, and F1).

execution accuracy, with smaller models exhibiting higher correlations between naturalness and performance. The presence of Least naturalness identifiers has the largest negative effect on schema linking. Additionally, while the performance difference between Regular and Low is visible, it is less impactful. So, modifying Least naturalness identifiers should be a higher priority than modifying Low naturalness identifiers.

5.1 Execution Accuracy

Execution Result Set Comparison. Execution accuracy is the standard measure of performance in most NL-to-SQL benchmarks [29, 61] where accuracy is determined using result set comparisons between gold and generated queries executed over one or more database instances. A drawback of existing methods is that strict set or bag comparisons risk increased false-negatives when a generated query includes additional fields that are not required, but do not render the result incorrect in terms of the natural language question [15, 62].

To reduce false negatives, the SNAIIS approach to execution accuracy evaluation adopts 2 aspects of relaxed execution matching as described in [15]; it accounts for: (1) The possibility that a predicted query may contain additional columns beyond those retrieved by a gold query; and (2) That unless specified in the NL question, tuples may appear in any order. To achieve this, we perform result set-superset comparisons to ensure that the predicted result set column set is a superset of the gold result set column set. That is, a generated query is considered incorrect if it does not contain *all* gold query columns; but it is not considered incorrect (at this stage) if it includes columns not present in the gold query result. A more detailed description of this approach is available in the appendix.

Manual Evaluation. Execution result set comparison cannot prove query correctness; so we rely on it only to rule out true negatives from further consideration. To validate correctness, the authors manually review generated queries that pass execution result set-superset comparison checks. We streamline this process by creating a Python-based manual validation user interface that makes the process of comparing gold and generated queries more user-friendly. Manual validation steps include ensuring the generated query answers the NL question, matches the gold query in terms of semantic structure, and does not contain semantically incorrect predicates, projections, or clauses.

Naturalness Effect on Execution Accuracy. Figure 8 shows execution accuracy for each LLM and naturalness level. There is a clear difference in overall performance between LLMs, most likely due to model size. We find that generally more natural database schemas yield more correct queries. Databases with more natural native schemas did not benefit from identifier renaming, though we observe that altering a schema to become less natural degrades accuracy in most cases. We find that for databases with Native schema combined naturalness scores less than 0.69, modifying the schema identifiers to increase naturalness improves execution accuracy.

Statistical Significance. The Kendall-Tau correlation between the naturalness of identifiers in a query and execution accuracy ranges from low ($\tau = 0.09, p < 0.0001$) for Gemini 1.5, to moderate ($\tau = 0.19, p < 0.0001$) for Phind-CodeLlama2 and CodeS. The most impactful relationship is between the presence of Least naturalness identifiers and performance, with Kendall-Tau correlations between the proportion of Least identifiers in a query and execution accuracy between $\tau = -.15$ and $\tau = -.22$ with $p < 0.0001$ for all models.

5.2 Schema Linking Evaluation

We make schema linking a “first class citizen” of our analysis, and study schema linking performance in queries irrespective of other aspects of correctness. Thus, we propose query-level and identifier-level schema linking measurements. We propose an approach similar to the Spider benchmark exact set matching system [61] in which we employ a schema linking-specific evaluation method using *recall* scoring of gold and generated query pairs. Other schema linking-focused research measure effects of schema linking improvements using ablation [6, 52, 58, 59]. In other cases, schema linking is described in post-hoc analysis of NL-to-SQL model performance, with schema linking accounting for roughly 30% of failures [12, 46].

Query-Level Linking Analysis. The set of all schema identifiers (table and column names) present in gold queries represents the minimum identifiers required to correctly answer an NL question. Our purpose-built ANTLR4-based [43] query parser extracts identifiers from gold and generated queries. With a set QI_g of identifiers present in the gold query and a set of identifiers QI_p present in the generated (or predicted) query, we calculate recall, as well as F1 and precision.

$$QueryRecall = \frac{|QI_g \cap QI_p|}{|QI_g|} \quad (1)$$

$$QueryPrecision = \frac{|QI_g \cap QI_p|}{|QI_p|} \quad (2)$$

$$QueryF1 = \frac{2(QueryRecall * QueryPrecision)}{QueryRecall + QueryPrecision} \quad (3)$$

We exclude 137 linking score calculations from analysis in situations where the predicted query contains invalid SQL that prevents query parsing and identifier extraction. We use recall as the primary measure for schema linking, as it does not penalize generated queries that contain extra identifiers that do not render an answer incorrect in our setting, such as cases when an arbitrary column is

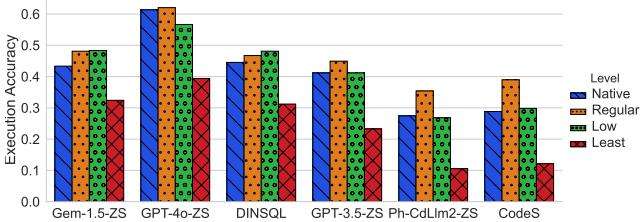


Figure 8: Execution accuracy (proportion of correct queries) by model. There is slight accuracy improvement from native schemas to schemas modified to regular naturalness. Accuracy drops significantly for schemas modified to low naturalness.

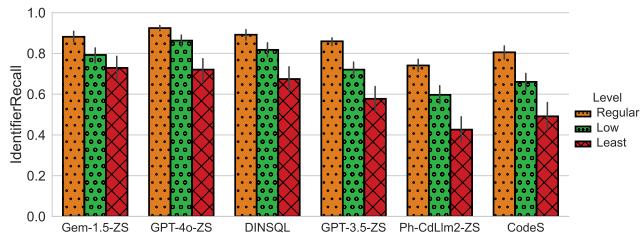


Figure 9: Native identifier recall scores by model and naturalness level. Error bars set with confidence interval of 0.95. For all models, identifiers in lower naturalness categories yield lower recall scores.

referenced in a count function. Charts and tables depicting F1 and precision scores are available in the appendix.

Identifier-Level Linking Analysis. For an identifier-focused (rather than query-focused) metric, we perform identifier-level linking analysis. We derive recall linking scores for each Native schema identifier I as follows. I_{match} is the count of instances when I is correctly present in a predicted query. I_{gold} is the count of gold queries that contain I .

$$\text{IdentifierRecall} = \frac{I_{match}}{I_{gold}} \quad (4)$$

Both DIN SQL and the CodeS complex NL-to-SQL workflows are sensitive to changes in naturalness, suggesting that these more complex workflows by themselves do not overcome schema naturalness effects. We also see that execution accuracy differences between the GPT-4o zero-shot prompting method and the DINSQL prompt chaining method suggest that applying more complex workflows to high-performing LLMs may be counterproductive for more recent SoTA LLMs.

Naturalness Effect on Schema Linking. Overall, we find that schema naturalness has a model-dependent and significant effect on schema linking performance with the highest correlations between QueryRecall and query naturalness occurring with the open-source CodeLlama and CodeS models, and the lowest (though still significant) correlations occurring with Google’s SoTA Gemini 1.5 Pro and OpenAI’s GPT-4o models. The more complex DIN SQL and

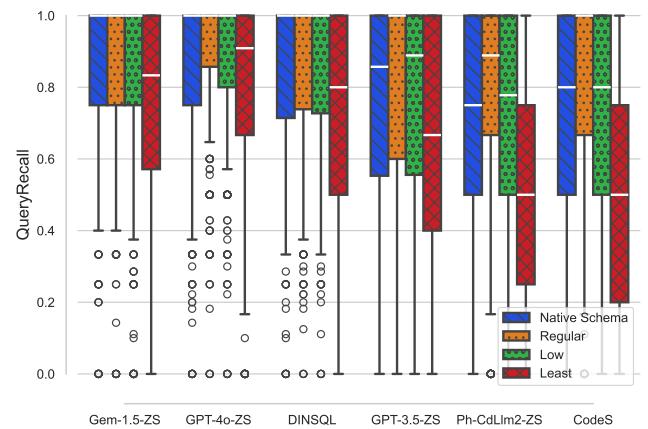


Figure 10: Schema linking performance across database schema naturalness levels. Generally yields equal or better performance for higher levels of naturalness, with open source models Phind-CodeLlama2 (Ph-CdLlm2-ZS) and CodeS as well as OpenAI’s GPT-3.5 (GPT-3.5-ZS) exhibiting higher sensitivity to changes in naturalness. Zero-shot prompting NL-to-SQL methods are denoted as (ZS).

CodeS workflow QueryRecall results are also significantly affected by naturalness level differences.

Figure 9 visualizes IdentifierRecall of Native identifiers in each naturalness level, and for each LLM. The chart indicates an observable difference in IdentifierRecall scores for each naturalness level, with IdentifierRecall increasing for higher naturalness levels. These results remain consistent relative to overall model performance across all 5 LLMs and various workflows.

Figure 10 illustrates QueryRecall across schema naturalness levels, and for each LLM. For GPT 3.5, Phind-CodeLlama2, and CodeS, we observe an improvement to QueryRecall when converting identifiers in a Native schema to Regular naturalness. This improvement did not manifest for Gemini and GPT-4o when observing the data in aggregate (i.e., between databases) due to their overall high performance relative to the other models, but improvements within databases of lower naturalness are still present (see Figure 11). The recall drop (approximately 20 percent decrease) associated with a modification from both Regular and Low to Least naturalness remains consistent across all LLMs.

Naturalness changes within specific SNAILS database schemas paints a clearer picture of the impact of naturalness. Figure 11 provides a drill-down view of the effect of schema modification on the PILB, SBOD, and NTSB schemas in terms of QueryRecall , and for each LLM and schema naturalness level. The center example (PILB) is a highly natural Native schema where schema naturalness modification would not be required. The leftmost example (NTSB) indicates linking performance improvement across all models for a native schema of lower naturalness converted to a higher naturalness schema, and presents a case where naturalness modification will improve NLI performance. The rightmost database (SBOD) represents a Least naturalness schema, and transformation from Native to Regular yields significant improvements for all models.

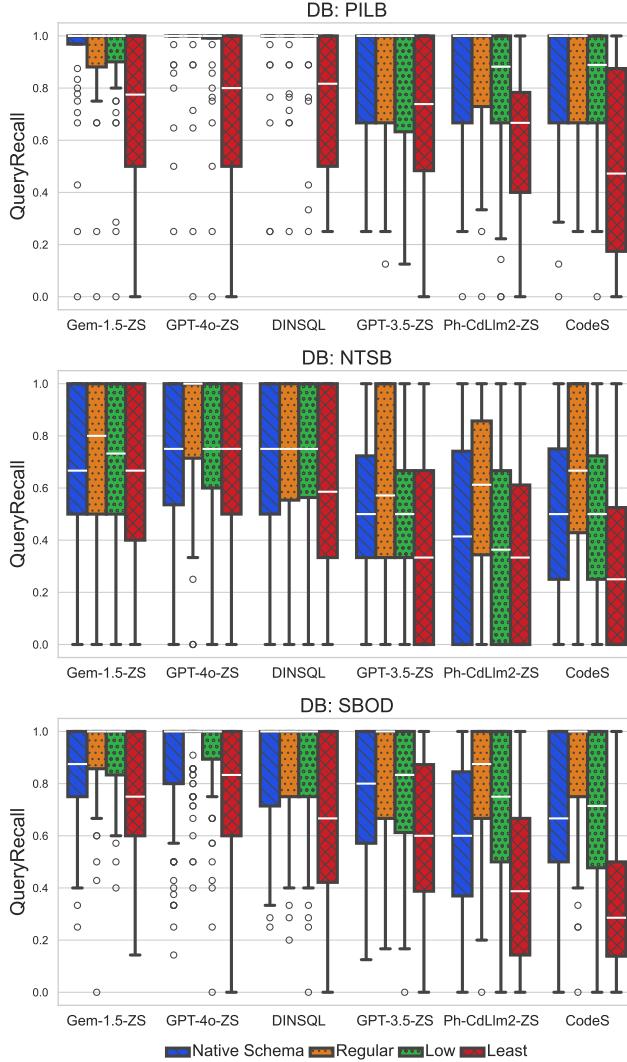


Figure 11: Schema linking performance (QueryRecall score) changes across 3 example databases’ native and virtual schemas. We selected these 3 examples to showcase the diversity of the databases in our collection. PILB Native is a more natural schema with 65 percent Regular, 22 percent Low, and 13 percent Least; NTSB Native contains 42 percent Regular, 34 percent Low, and 24 percent Least; and SBOD Native is the lowest naturalness schema with 24 percent Regular, 49 percent Low, and 27 percent Least.

In all cases, we see that reducing naturalness to the Least level consistently degrades *QueryRecall*.

Statistical Significance. Kendall-Tau correlations between the proportion of Least identifiers and *QueryRecall* range from $\tau = -0.16$ (Gemini) to $\tau = -0.28$ (Phind-CodeLlama2), with $P < 0.001$ for all models. Both Regular and Low identifier proportions are significantly correlated with improved outcomes in terms of *QueryRecall*. Identifiers with Regular naturalness show the highest positive

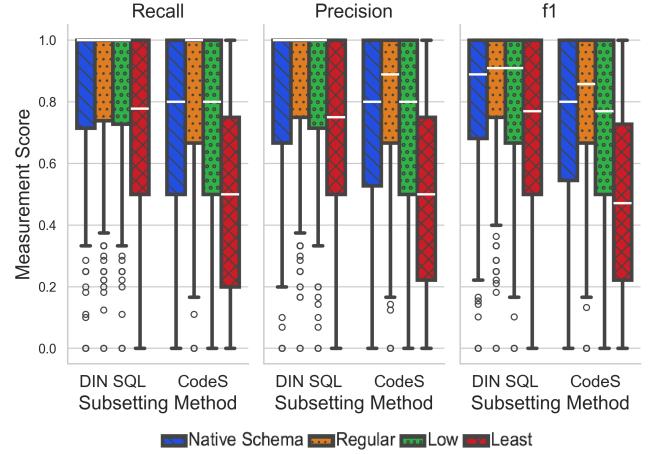


Figure 12: Schema subsetting performance, measured with recall, precision, and f1 score, varies by naturalness levels for both DIN SQL and CodeS. Measurement Score is Recall, Precision, or f1 respectively.

Kendall-Tau correlations ranging from $\tau = 0.07$ (Gemini) to $\tau = 0.20$ (Phind-CodeLlama2). Low naturalness identifier proportions correlate positively, but to a lesser degree, with Kendall-Tau values ranging from $\tau = 0.05$ (Phind-CodeLlama2) to $\tau = 0.07$ (Gemini).

Naturalness Effects on Schema Subsetting. We measure the schema subsetting (also known as schema filtering, or table retrieval) in terms of recall, precision, and f1 score, and present the results in Figure 12. We find that for the CodeS finetuned classifier approach, schema naturalness level differences result in observable differences in f1. For the DIN SQL LLM-based approach, naturalness effects are less pronounced, though still present, particularly for Least level schemas.

Performance Over Modified Spider Schemas. Figure 13 shows that with the SNAILS schema renaming artifacts applied to the Spider NL-to-SQL benchmark dev dataset [61], naturalness effects are the most significant between Low and Least levels of naturalness. Performance differences across naturalness levels for the highly natural Spider schemas resemble performance over similarly-natural schemas in the SNAILS collection.

Additional Charts and Figures. The appendix also provides additional fine-grained results: a more detailed tabular breakdown of execution accuracy by schema and LLM; Precision- and F1-based results; token ratio correlations; and more granular *QueryRecall* correlations and box plots;

6 DISCUSSION AND LIMITATIONS

The ability to assess the naturalness of existing schemas can inform the feasibility of “hooking up” an NL query interface to an existing database. We believe that practitioners who are considering the integration of an LLM into their database interaction workflows would benefit from naturalness-focused schema analysis a key step in their integration process.

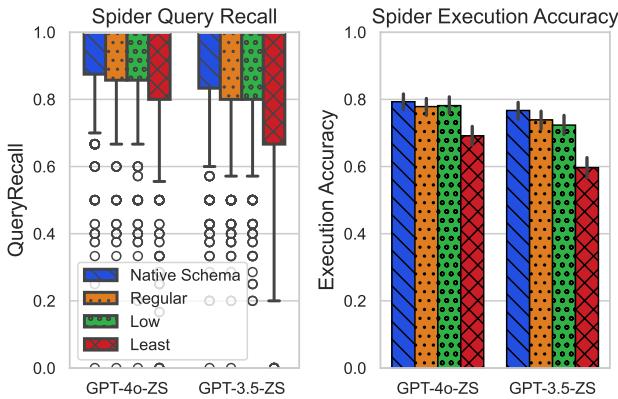


Figure 13: QueryRecall and Execution Accuracy differences over the Spider [61] dev set modified using SAILS renaming artifacts.

Other Naming Patterns in Real-World Schemas. To examine naming practices in the real-world, we classified the identifiers of SchemaPile dataset [11] with our CANINE-based classifier, and evaluated the identifiers for other LLM-unfriendly patterns. We observe that whitespace characters within schema identifiers contributes to identifier mutation during inference. That is, rather than encasing whitespace-containing identifier with brackets or quotes, the LLM hallucinates the identifier into snake or camel case format. We find that in the SchemaPile collection, though whitespace is uncommon (less than 1 percent for both tables and columns), it appears in 808 columns and 63 tables, and is comparable to the proportions in the SAILS dataset.

Another naming practice that yields disproportionate failures with some LLMs is the presence of the word *table* in the identifier name. In these instances, we find that the LLM tends to drop the word *table* from the name (e.g., `table_employee` becomes `employee`). There are over 700 identifiers (less than 1 percent of all identifiers) in the SchemaPile collection that employ this naming pattern.

These observations suggest that although these naming patterns are not necessarily a common occurrence in many real-world schema designs, they do appear. We suggest that practitioners would benefit from assessing the naming patterns of their database schemas.

Variations in LLM Sensitivity to Naturalness. There are many LLMs to select from for NLIDBs, and we can see even within the select 5 models in our work large variations in NL-to-SQL performance as well as the degree of sensitivity to schema naturalness. The Google Gemini and GPT-4o models demonstrate the highest overall performance, as well as the lowest sensitivity to naturalness differences between Regular and Low levels. Without access to the underlying model architectures and weights, it remains as a black box in our research, and we can merely speculate the reasons why it is not as affected by naturalness as the other 3 models in our study. Generally, we observe that these models have an overall higher performance, and are less prone to linking errors such as selecting

the incorrect identifier from the schema knowledge representation or committing a typo-like hallucination.

Though selecting the most performant model would seem to be an obvious course of action, competing factors such as an organization’s policies, budget, or existing vendor contracts, may require the selection of a model that is more sensitive to schema naturalness differences. Thus, we believe that naturalness-aware NLI integration will remain important for at least the practitioners who use LLMs other than Gemini in the set that we have studied.

Modifying Existing Schemas. For already-existing schemas, renaming identifiers is generally a non-trivial effort, particularly for those databases for which documentation has been published and application interfaces have been integrated. Schema modifications may not be necessary (or helpful), if a schema is already classified as highly natural. DBAs should assess current naturalness levels prior to committing to naming modifications. At a minimum, we recommend that any Least identifier be modified to a Regular naturalness level and, if feasible, Low identifiers as well. If renaming a less natural schema’s identifiers is not feasible due to integration constraints, we suggest one of two approaches: 1) adopting a naturalness-as-a-view strategy by mapping Native identifiers to Regular naturalness identifiers using SQL views, or 2) a middleware approach that modifies schema knowledge in LLM prompts generated queries prior to execution on the database. We sketch a rough design of **both options** in the appendix.

We demonstrate a natural schema view proof of concept with our SAILS database collection and their MS SQL Server instances. For each table and column in the collection’s database schemas, we map the Native table or column to its Regular counterpart in the naturalness modified identifier dataset using SQL view creation DDL and a `db_nl` schema. This enables schema information retrieval for LLM-based NL-to-SQL prompting without prompt or generated query modification while still retaining the underlying Native schema naming patterns required for existing integrations.

In lieu of schema modification, practitioners may elect to employ prompting techniques that augment schema representations with additional metadata or value samples. While these methods may improve schema linking performance in some contexts [34], they greatly increase schema representations on a per-identifier basis. Thus, the cost to do so is high in terms of token efficiency, latency, and implementation complexity, especially for very large schemas.

Designing New Schemas. For new schema development, our results show that making schema identifiers more natural from the start can make databases work better with LLMs. Specifically, database designers should try to avoid Least naturalness identifiers and would likely also benefit from limiting Low naturalness identifiers. Database practitioners can evaluate the naturalness of identifiers using the identifier naturalness classification techniques and model artifacts described in this paper and released publicly by us as part of the SAILS collection.

Limitations. LLM research is advancing rapidly, and the LLMs represented in this paper may get superseded by newer versions or newer models (e.g., DBRX [56], Arctic [57]). But it does not negate our work’s core value—the first in-depth characterization of how schema naturalness affects LLM-based NL-to-SQL—and our new

labeled datasets, AI artifacts, and benchmarking framework can be used for future LLMs too. We leave it to future work to also include such very recent LLMs for further benchmark analyses.

We recognize that the correlation statistics indicate a moderate (in some cases only a weak) correlation between naturalness and *IdentifierRecall*. This suggests that other undiscovered factors also influence linking performance; and further research may reveal additional schema- and language-related correlations.

Our selection of 9 database schemas is of course not fully representative of *all* types of schemas available in the real-world. The SAILS collection will benefit from continued growth in terms of both databases and NL-SQL pairs. We hope our open source datasets and artifacts can be built upon by the database and NLP communities to keep improving LLM-based NL-to-SQL.

Future Work. In addition to extending the SAILS benchmark artifacts to include additional datasets and artifacts, we identify several NLP+DB directions for future work. First, we wish to ask why and how exactly do different naturalness levels alter schema linking performance so much? Is it due to the tokenization and embedding mechanics? If so, where in the latent space do these altered tokens end up, and how do the encoders make use of them? Second, why do the different foundational LLMs behave so differently? Is it related to their architectures, tokenization, (pre)training data, post-training finetuning process, or some other factors? We believe these open questions have the potential to lead to several interesting new lines of research at the DB and NLP intersection.

7 RELATED WORK

Ontology Mapping. Schema modifications and intermediate representations to enhance performance in a specific context extend beyond NL-to-SQL applications. Mapping relational database schemas to ontologies is an approach used to improve schema-to-schema integration and web application application-database interfaces [60]. This improves the semantic description of underlying data, which is often a desirable feature in web applications that interact within the semantic web [21]. While ontological mapping of a relational database can improve performance in this context; we see less evidence that such an approach is useful or necessary in NL-to-SQL applications, though this may serve as a compelling opportunity for future research.

NL-to-SQL Benchmarks. *Spider* [61], soon to be superseded by a more challenging benchmark for the LLM era, was a popular NL-to-SQL benchmark that still offers a publicly-available dataset consisting of 166 multi-table databases and 1,034 NL questions and gold queries over the databases in a development dataset. *Spider-Syn* [16] and *Spider-Realistic* [16] are extensions of the Spider benchmark that perform NL question synonym replacement to reduce the occurrences of lexical matching between NL question keywords and schema identifiers. *BIRD* [29] is an emergent benchmark containing 95 large databases over 37 domains that seeks to better replicate real-world databases in order to better challenge highly capable LLM-based NL-to-SQL systems. While Spider and its variants as well as BIRD intend to better-replicate real-world database designs, our naturalness-focused analysis indicates that their schema identifiers are more natural than those we encountered

in our real-world database selection process (see the statistics in Figure 3). Additionally, Spider and BIRD both evaluate performance using either exact set matching or execution result set comparison while we use the more pragmatic set-superset matching as proposed in [15] and schema linking-specific recall metrics.

Archerfish [15] is a benchmarking framework that relaxes execution matching and accounts for semantic ambiguity in NL questions by allowing for multiple correct answers derived from candidate key analysis. This framework relies on the binary “correct, or not” evaluation approach common to other benchmarks, whereas in addition to relaxed execution matching, SAILS evaluates target schema linking performance via query identifier recall. Overall, we find that our benchmark and findings complement this existing and ongoing research by enhancing our ability to target specific schema-related aspects of NL-to-SQL performance in future NLI development.

Impacts of Schema on NL-to-SQL Performance. *Spider-Syn* [16] demonstrates degraded NL-to-SQL performance of language models trained for NL-to-SQL tasks when the occurrence of lexical matching between NL questions and schema identifiers is reduced. This approach differs from our experiments in that it evaluates a LM specifically trained on NL-to-SQL tasks using the Spider training set as opposed to the more general-purpose foundational LLMs evaluated in this work. They also make no apparent attempt to reduce the naturalness of database schema identifiers.

Semantics-preserving schema transformation is a design feature of MT-teql [31], an NL-to-SQL evaluation framework that modifies natural language utterances and schema properties to stress LM robustness. MT-teql provides a holistic view of the effect of NL utterance variances and schema design on LM performance. However, it does not address the question of schema identifier naturalness, nor does it make modifications to schema elements that are necessary for answer generation.

Some recent work has examined the effects of schema ambiguity, where semantically different tables or columns have identical or synonymous names. Schema ambiguity, where a schema contains one or more semantically similar pairs of elements, degrades semantic parsing (i.e., NL-to-SQL) performance by recalling undesired tables or columns in response to a NL question that contains patterns or keywords that align with more than one schema element in the latent space [41]. Documentation, combined with agent-based column selection, can improve Text-to-SQL performance in the presence of data and schema ambiguity [22]. Though we did not focus on ambiguity in our work, [identifier naturalness and ambiguity are complementary efforts that provide a potential future direction for the expansion of the SAILS benchmark artifacts](#).

ACKNOWLEDGMENTS

We thank the members of the UCSD Database Lab and Jingbo Shang for their help reviewing this work. This work was supported in part by the U.S. Army Advanced Civil Schooling program, and gifts from VMWare and CloudBank. LLMs were used for Latex table and diagram formatting.

REFERENCES

- [1] [n.d.]. NPS IRMA Portal. <https://irma.nps.gov/Portal/>. Accessed: April 2023.

- [2] [n.d.] SAP TABLES. <https://sap.erpref.com/>. Accessed: June 2023.
- [3] 2022. grammars-v4. <https://github.com/antlr/grammars-v4>.
- [4] 2022. Report Card Database 2021-22. <https://data.nysed.gov/files/essa/21-22/SRC2022.zip>. Accessed: May 2023.
- [5] Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sumita Dev, Jacob Devlin, Mark Diaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeonjae Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Miggiani, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee Shelby, Ambrose Sloane, Daniel Smilkov, David R. So, Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhua Wu, Kelvin Xu, Yunhan Xu, Linting Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav Petrov, and Yonghui Wu. 2023. PaLM 2 Technical Report. arXiv:2305.10403 [cs.CL]
- [6] Ruisheng Cao, Lu Chen, Zhi Chen, Yanbin Zhao, Su Zhu, and Kai Yu. 2021. LGESQL: Line Graph Enhanced Text-to-SQL Model with Mixed Local and Non-Local Relations. 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)*. Association for Computational Linguistics, Online, 2541–2555. <https://doi.org/10.18653/v1/2021.acl-long.198>
- [7] Pankaj Kumar Choudhary. 2022. Naming Conventions in SQL. <https://www.c-sharpcorner.com/UploadFile/f0b2ed/what-is-naming-convention/>. Last accessed on 2024-01-01.
- [8] Jonathan H. Clark, Dan Garrette, Iulia Turc, and John Wieting. 2022. Canine: Pre-training an Efficient Tokenization-Free Encoder for Language Representation. *Transactions of the Association for Computational Linguistics* 10 (2022), 73–91. https://doi.org/10.1162/tacl_a_00448
- [9] Robert Cook. 2016. Field Data for Assateague Island National Seashore Amphibian and Reptile Inventory. <https://irma.nps.gov/DataStore/Reference/Profile/2236826>. Accessed: April 2023.
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. arXiv:1810.04805 [cs.CL]
- [11] Till Doehmen, Radu Geacu, Madelon Hulsebos, and Sebastian Schelter. 2024. SchemaPile: A Large Collection of Relational Database Schemas. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*.
- [12] Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao, Yunjun Gao, lu Chen, Jinshi Lin, and Dongfang Lou. 2023. C3: Zero-shot Text-to-SQL with ChatGPT. arXiv:2307.07306 [cs.CL]
- [13] Thomas Evans. 2015. Great Smoky Mountains All Taxa Biodiversity Inventory (ATBI) Plot Vegetation Monitoring Database. <https://irma.nps.gov/DataStore/Reference/Profile/2221324>. Accessed: April 2023.
- [14] Catherine Finegan-Dollak, Jonathan K. Kummerfeld, Li Zhang, Karthik Ramamathan, Sesh Sadasivam, Rui Zhang, and Dragomir Radev. 2018. Improving Text-to-SQL Evaluation Methodology. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Melbourne, Australia, 351–360. <https://doi.org/10.18653/v1/P18-1033>
- [15] Avrilia Floratou, Fotis Psallidas, Fuhseng Zhao, Shaheed Deep, Gunther Hagleitner, Wangda Tan, Joyce Cahoon, Rana Aloabi, Jordan Henkel, Abhik Singla, Alex Van Grootel, Brandon Chow, Kai Deng, Katherine Lin, Marcos Campos, Venkatesh Emani, Vivek Pandit, Victor Shnayder, Wenjing Wang, and Carlo Curino. 2024. NL2SQL is a solved problem... Not!. In *Proceedings of the CIDRDB 2024 Conference*. <https://www.cidrdb.org/cidr2024/papers/p74-floratou.pdf>
- [16] Yujian Gan, Xinyu Chen, Qiuping Huang, Matthew Purver, John R. Woodward, Jinxia Xie, and Pengsheng Huang. 2021. Towards Robustness of Text-to-SQL Models against Synonym Substitution. 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)*. Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (Eds.). Association for Computational Linguistics, Online, 2505–2515. <https://doi.org/10.18653/v1/2021.acl-long.195>
- [17] Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, Yichen Qian, Bolin Ding, and Jingren Zhou. 2024. Text-to-SQL Empowered by Large Language Models: A Benchmark Evaluation. , 14 pages. <https://doi.org/10.14778/3641204.3641221>
- [18] General Services Administration. 2023. Security Policy for Generative Artificial Intelligence (AI) Large Language Models (LLMs). <https://www.gsa.gov/directives-library/security-policy-for-generative-artificial-intelligence-ai-large-language-models-llms>. Last accessed on 2024-05-28.
- [19] Frederic Piesschaert Gert Van Spaendonk, Jo Loos. [n.d.] Database naming conventions. https://inbo.github.io/tutorials/tutorials/database_conventions/. Last accessed on 2024-01-01.
- [20] Sree Hari Krishnan Parthasarathi, Lu Zeng, and Dilek Hakkani-Tür. 2023. Conversational Text-to-SQL: An Odyssey into State-of-the-Art and Challenges Ahead. In *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1–5. <https://doi.org/10.1109/ICASSP49357.2023.10096170>
- [21] Mohamed A. G. Hazber, Ruixuan Li, Yuxi Zhang, and Guandong Xu. 2015. An Approach for Mapping Relational Database into Ontology. In *2015 12th Web Information System and Application Conference (WISA)*. 120–125. <https://doi.org/10.1109/WISA.2015.25>
- [22] Zezhou Huang, Pavan Kalyan Damalapati, and Eugene Wu. 2023. Data Ambiguity Strikes Back: How Documentation Improves GPT’s Text-to-SQL. In *NeurIPS 2023 Second Table Representation Learning Workshop*. <https://openreview.net/forum?id=FhKTuIRTd>
- [23] Klamath Inventory and Monitoring Network. 2021. Exotic and Invasive Plants Monitoring Database. <https://irma.nps.gov/DataStore/Reference/Profile/2288667>. Accessed: April 2023.
- [24] Seth Judge and Kevin Kozar. 2023. Pacific Island Network Landbird Monitoring Dataset. <https://irma.nps.gov/DataStore/Reference/Profile/2300107>. <https://doi.org/10.57830/2300107> Accessed: April 2023.
- [25] Tobias Kuhn. 2014. A Survey and Classification of Controlled Natural Languages. *Computational Linguistics* 40, 1 (03 2014), 121–170. https://doi.org/10.1162/COLI_a_00168 arXiv:https://direct.mit.edu/coli/article-pdf/40/1/121/1812691/coli_a_00168.pdf
- [26] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- [27] Haoyang Li, Jing Zhang, Cuiping Li, and Hong Chen. 2023. RESDSL: decoupling schema linking and skeleton parsing for text-to-SQL. , 9 pages. <https://doi.org/10.1609/aaai.v37i11.26535>
- [28] Haoyang Li, Jing Zhang, Hanbing Liu, Ju Fan, Xiaokang Zhang, Jun Zhu, Renjie Wei, Hongyan Pan, Cuiping Li, and Hong Chen. 2024. CodeS: Towards Building Open-source Language Models for Text-to-SQL. *Proc. ACM Manag. Data* 2, 3, Article 127 (May 2024), 28 pages. <https://doi.org/10.1145/3654930>
- [29] Jinyang Li, Bin Yuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, Xuanhe Zhou, Chenhao Ma, Guoliang Li, Kevin C.C. Chang, Fei Huang, Reynold Cheng, and Yongbin Li. 2024. Can LLM already serve as a database interface? a big bench for large-scale database grounded text-to-SQLs. , 28 pages.
- [30] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Deniz Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebazé, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblikulov, Zhiruo Wang, Rudra Murthy, Jason Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. 2023. StarCoder: may the source be with you! (2023). arXiv:2305.06161 [cs.CL]
- [31] Pingchuan Ma and Shuai Wang. 2021. MT-Teql: Evaluating and Augmenting Neural NLIDB on Real-World Linguistic and Schema Variations. *Proc. VLDB Endow.* 15, 3 (nov 2021), 569–582. <https://doi.org/10.14778/3494124.3494139>
- [32] Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veysel, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. 2023. Recent Advances in Natural Language Processing via Large Pre-Trained Language Models: A Survey. *ACM Comput. Surv.* 56, 2, Article 30 (sep 2023), 40 pages. <https://doi.org/10.1145/3605943>
- [33] Ian Muirhead. 2021. Northern Great Plains Fire Management: FFI Database. <https://irma.nps.gov/DataStore/Reference/Profile/2297267>. Accessed: April 2023.
- [34] Linyong Nan, Yilun Zhao, Weijin Zou, Narutatsu Ri, Jaesung Tae, Ellen Zhang, Arman Cohan, and Dragomir Radev. 2023. Enhancing Few-shot Text-to-SQL Capabilities of Large Language Models: A Study on Prompt Design Strategies. arXiv:2305.12586 [cs.CL] <https://arxiv.org/abs/2305.12586>
- [35] National Center for Statistics and Analysis. 2022. Overview of the 2021 Crash Investigation Sampling System. <https://crashstats.nhtsa.dot.gov/Api/Public/>

- ViewPublication/813397 Traffic Safety Facts Research Note. Report No. DOT HS 813 397.
- [36] OpenAI. [2022]. ChatGPT: Optimizing Language Models for Dialogue. <https://openai.com/blog/chatgpt/>.
- [37] OpenAI. 2023. OpenAI API Documentation. <https://platform.openai.com/docs/guides/gpt>. Last accessed on 2023-10-30.
- [38] OpenAI. 2023. OpenAI Tokenizer. <https://github.com/openai/tiktoken>. Last accessed on 2023-10-30.
- [39] Oracle. [n.d.]. Database Object Names and Qualifiers. <https://docs.oracle.com/en/database/oracle/oracle-database/19/sqlrf/Database-Object-Names-and-Qualifiers.html>. Last accessed on 2024-01-01.
- [40] Oracle. [n.d.]. Table Naming Standards and Conventions. https://docs.oracle.com/cd/E92917_01/PDF/8.1.x/common/HTML/DM_Naming/2_Table_and_Column_Naming_Standards.htm. Last accessed on 2023-11-07.
- [41] Simone Papicchio, Paolo Papotti, and Luca Cagliero. 2024. Evaluating Ambiguous Questions in Semantic Parsing. In *2024 IEEE 40th International Conference on Data Engineering Workshops (ICDEW)*. 338–342. <https://doi.org/10.1109/ICDEW61823.2024.00050>
- [42] Simone Papicchio, Paolo Papotti, and Luca Cagliero. 2024. QATCH: benchmarking SQL-centric tasks with table representation learning models on your data. In *Proceedings of the 37th International Conference on Neural Information Processing Systems* (New Orleans, LA, USA) (NIPS '23). Curran Associates Inc., Red Hook, NY, USA, Article 1348, 20 pages.
- [43] Terence Parr, Sam Harwell, and Kathleen Fisher. 2014. Adaptive LL(*) Parsing: The Power of Dynamic Analysis. In *Proceedings of the 2014 ACM International Conference on Object Oriented Programming Systems Languages and Applications* (Portland, Oregon, USA) (OOPSLA '14). Association for Computing Machinery, New York, NY, USA, 579–598. <https://doi.org/10.1145/2660193.2660202>
- [44] Phind. 2023. Phind-CodeLlama-34B-v2. <https://huggingface.co/Phind/Phind-CodeLlama-34B-v2>.
- [45] Marie-Laurence Poujois. 2021. Localized Demo Databases Now Available for SAP Business One 10.0 FP 2011. <https://blogs.sap.com/2021/01/29/localized-demo-databases-now-available-for-sap-business-one-10-0-fp-2011/>. Accessed: April 2023.
- [46] Mohammadreza Pourreza and Davood Rafie. 2024. DIN-SQL: decomposed in-context learning of text-to-SQL with self-correction. In *Proceedings of the 37th International Conference on Neural Information Processing Systems* (New Orleans, LA, USA) (NIPS '23). Curran Associates Inc., Red Hook, NY, USA, Article 1577, 10 pages.
- [47] G. A. Radja, E.-Y. Noh, and F. Zhang. 2022. Crash Investigation Sampling System 2021 analytical user's manual. Accessed: April 2023.
- [48] Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoaqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémie Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Cantor Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. Code Llama: Open Foundation Models for Code. arXiv:2308.12950 [cs.CL]
- [49] Tushar Sharma, Marios Fragnoulis, Stamatia Rizou, Magiel Bruntink, and Diomidis Spinellis. 2018. Smelly Relations: Measuring and Understanding Database Schema Quality. In *Proceedings of the 40th International Conference on Software Engineering: Software Engineering in Practice* (Gothenburg, Sweden) (ICSE-SEIP '18). Association for Computing Machinery, New York, NY, USA, 55–64. <https://doi.org/10.1145/3183519.3183529>
- [50] StackOverflow. [n.d.]. Database, Table and Column Naming Conventions? <https://stackoverflow.com/questions/7662/database-table-and-column-naming-conventions>. Last accessed on 2024-01-01.
- [51] Charles Stefanic. 2021. Wildlife Observations Database: Craters of the Moon National Monument and Preserve 1921–2021. <https://irma.nps.gov/DataStore/Reference/Profile/2192964>. Accessed: April 2023.
- [52] Alane Laughlin Suhr, Kenton Lee, Ming-Wei Chang, and Pete Shaw. 2020. Exploring Unexplored Generalization Challenges for Cross-Database Semantic Parsing. In *ACL 2020*.
- [53] Gemini Team. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv:2403.05530 [cs.CL]
- [54] Gemini Team. 2024. Gemini: A Family of Highly Capable Multimodal Models. arXiv:2312.11805 [cs.CL]
- [55] Gemma Team. 2024. Gemma: Open Models Based on Gemini Research and Technology. arXiv:2403.08295 [cs.CL]
- [56] The Mosaic Research Team. 2024. Introducing DBRX: A New State-of-the-Art Open LLM. <https://www.databricks.com/blog/introducing-dbrx-new-state-of-the-art-open-llm>.
- [57] The Snowflake Research Team. 2024. Snowflake Arctic: The Best LLM for Enterprise AI - Efficiently Intelligent, Truly Open. <https://www.snowflake.com/blog/arctic-open-efficient-foundation-language-models-snowflake>.
- [58] Bailin Wang, Richard Shin, Xiaodong Liu, Alex Polozov, and Matthew Richardson. 2020. RAT-SQL: Relation-Aware Schema Encoding and Linking for Text-to-SQL Parsers. In *ACL 2020*. <https://www.microsoft.com/en-us/research/publication/rat-sql-relation-aware-schema-encoding-and-linking-for-text-to-sql-parsers/>
- [59] Lihan Wang, Bowen Qin, Binyuan Hui, Bowen Li, Min Yang, Bailin Wang, Binhua Li, Jian Sun, Fei Huang, Luo Si, and Yongbin Li. 2022. Proton: Probing Schema Linking Information from Pre-Trained Language Models for Text-to-SQL Parsing. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (Washington DC, USA) (KDD '22). Association for Computing Machinery, New York, NY, USA, 1889–1898. <https://doi.org/10.1145/3534678.3539305>
- [60] Zhuoming Xu, Shichao Zhang, and Yisheng Dong. 2006. Mapping between Relational Database Schema and OWL Ontology for Deep Annotation. In *2006 IEEE/WIC/ACM International Conference on Web Intelligence (WI 2006 Main Conference Proceedings) (WI'06)*. 548–552. <https://doi.org/10.1109/WI.2006.114>
- [61] Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxi Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2018. Spider: A Large-Scale Human-Labeled Dataset for Complex and Cross-Domain Semantic Parsing and Text-to-SQL Task. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium.
- [62] Lu Zeng, Sree Hari Krishnan Parthasarathi, and Dilek Hakkani-Tur. 2023. N-Best Hypotheses Reranking for Text-to-SQL Systems. In *2022 IEEE Spoken Language Technology Workshop (SLT)*. 663–670. <https://doi.org/10.1109/SLT54892.2023.10023434>
- [63] Jian Zhang, Zhengyuan Shen, Balasubramanian Srinivasan, Shen Wang, Huzeifa Rangwala, and George Karypis. 2023. NameGuess: Column Name Expansion for Tabular Data. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 13276–13290. <https://doi.org/10.18653/v1/2023.emnlp-main.820>

A SOURCE DATABASES AND QUERIES IN SAILS

The SAILS database collection (Artifact 1) contains 9 real-world schemas sourced from several different organizations and domains. In this section, we describe the technical details of each schema and its associated NL question-SQL query pairs. The bar charts in each data source section visually portray the clause compositions of the gold queries in Artifact 6—the NL question-SQL query pairs.

A.1 Data Sources

A.1.1 Field Data for Assateague Island National Seashore Amphibian and Reptile Inventory (ASIS).

Data Description. The ASIS database [9] is sourced from the National Parks Service (NPS) Irma portal [1] and contains scientific observation data of wildlife in the Assateague Island National Seashore preserve.

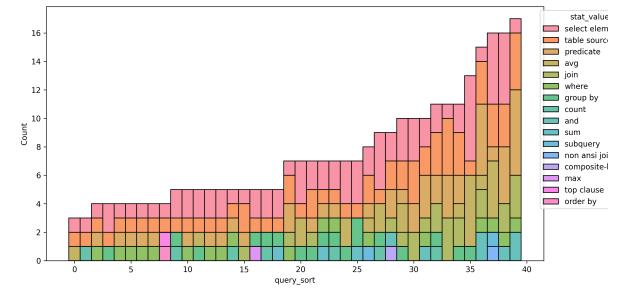


Figure 14: Gold query clause composition - ASIS database

ASIS Database Technical Details.

- Data source format: Microsoft Access

- Migration method: SQL Server Migration Assistant
- Table count: 36
- Column count: 245
- Mean columns per table: 6.125
- NL Questions: 40
- Combined naturalness level: 0.77

A.1.2 Great Smoky Mountains All Taxa Biodiversity Inventory (ATBI) Plot Vegetation Monitoring Database.

Data Description. The ATBI database [13] contains scientific observations of vegetation within the Great Smoky Mountains national park.

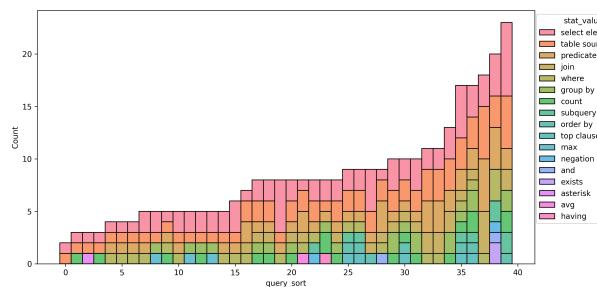


Figure 15: Gold query clause composition - ATBI database

ATBI Database Technical Details.

- Data source format: Microsoft Access
- Migration method: SQL Server Migration Assistant
- Table count: 28
- Column count: 192
- Mean columns per table: 6.857
- NL Questions: 40
- Combined naturalness level: 0.70

A.1.3 Klamath Inventory and Monitoring Network (KIS).

Data Description. The Klamath Invasive Species (KIS) database [23] contains scientific observations of exotic and invasive plants observed in Klamath Falls, Oregon.

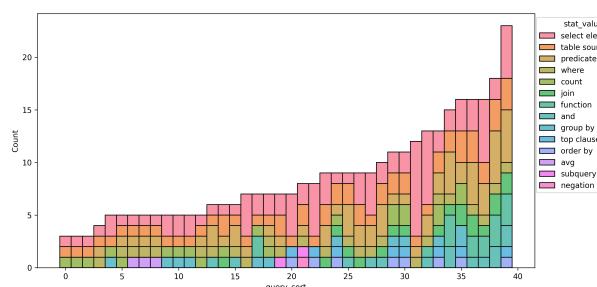


Figure 16: Gold query clause composition - KIS database

KIS Database Technical Details.

- Data source format: Microsoft Access
- Migration method: SQL Server Migration Assistant
- Table count: 18
- Column count: 157
- Mean columns per table: 8.72
- NL Questions: 40
- Combined naturalness level: 0.79

A.1.4 Pacific Island Network Landbird Monitoring Dataset.

Data Description. Pacific island landbirds (PILB) database [24] contains scientific observations of bird observations in various Pacific islands within the US states and territories.

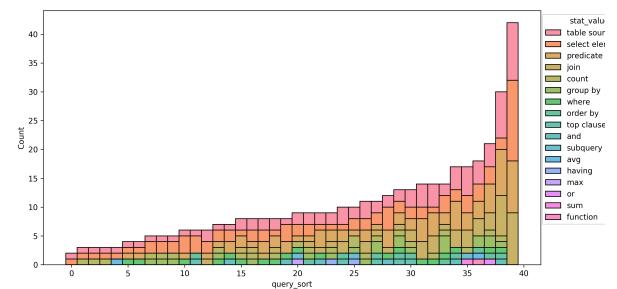


Figure 17: Gold query clause composition - PILB database

PILB Database Technical Details.

- Data source format: Microsoft Access
- Migration method: SQL Server Migration Assistant
- Table count: 21
- Column count: 196
- Mean columns per table: 9.33
- NL Questions: 40
- Combined naturalness level: 0.75

A.1.5 Wildlife Observations Database: Craters of the Moon National Monument and Preserve 1921–2021.

Data Description. The Craters Wildlife Observation (CWO) database [51] contains observations of wildlife spotted at the Craters of the Moon national monument and preserve. It is the smallest and most natural database in the benchmark data set.

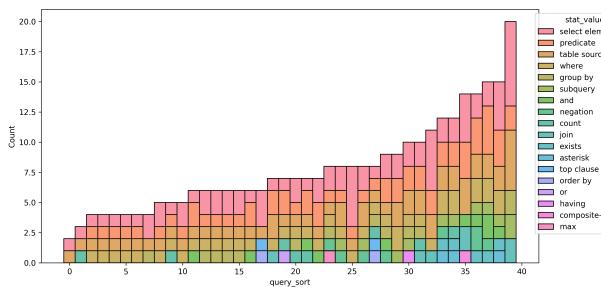


Figure 18: Gold query clause composition - CWO database

CWO Database Technical Details.

- Data source format: Microsoft Access
- Migration method: SQL Server Migration Assistant
- Table count: 13
- Column count: 71
- Mean columns per table: 5.461
- NL Questions: 40
- Combined naturalness level: 0.84

A.1.6 Northern Great Plains Fire Management: FFI Database.

Data Description. The NPFM database [33] contains observations of various overstory and other flora within the Northern Plains region of the National Parks Service.

Code-Bison Evidence of Familiarity. With this dataset, we observe some indications of exposure to the Code-Bison language model. We note that we no longer report performance results of Code-Bison inference in our main report.

When prompted with the NL question "How many overstory's have a codominant canopy position?", it generated the query:

```
SELECT COUNT(*)
FROM tbl_Overstory
WHERE CanPos = 2;
```

Which is a correct reference to the canopy position (CanPos) lookup code of 2, which corresponds to the codominant canopy position. The LLM was not provided code lookup information within the prompt, which suggests that some reference to the NPFM schema was included in its training data.

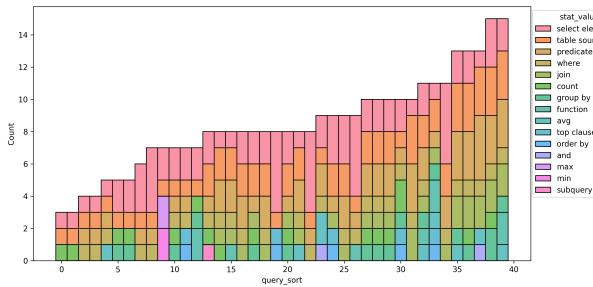


Figure 19: Gold query clause composition - NPFM database

NPFM Database Technical Details.

- Data source format: Microsoft Access
- Migration method: SQL Server Migration Assistant
- Table count: 27
- Column count: 190
- Mean columns per table: 7.037
- NL Questions: 40
- Combined naturalness level: 0.70

A.1.7 2021 Crash Investigation Sampling System.

Data Description. The crash investigation sampling system [47] is sourced from the National Transportation Safety Board, and referred to as NTSB in this paper. It contains sampled data of crash and vehicle statistics from 2021. The data is organized such that composite key joins are required for most multi-relation queries.

Additional Implementation Details. This is the only database in our collection that required deliberate migration from a non-database format to the target MS SQL Server environment. We acquired the data in .csv form, with a single .csv per table. Analysis of the files confirmed that although not in database form, the data was relational in nature, and migration involved SQL-based ingestion from .csv files into the target schema. The .sql scripts used to generate the database schema and insert table values is available in the project repository.

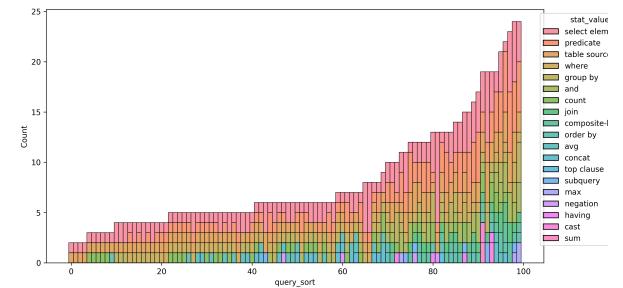


Figure 20: Gold query clause composition - NTSB database

NTSB Database Technical Details.

- Data source format: Comma Separated Value (CSV) files
- Migration method: SQL database creation and Python-based ETL scripting
- Table count: 40
- Column count: 1,611
- Mean columns per table: 40.275
- NL Questions: 100
- Combined naturalness level: 0.59

A.1.8 New York State Education Department Report Card Database 2021-22.

Data Description. The NYSED database [4] is sourced from the New York State Education Department. It contains standardized testing and demographic data for all public elementary, middle, and high schools in New York State.

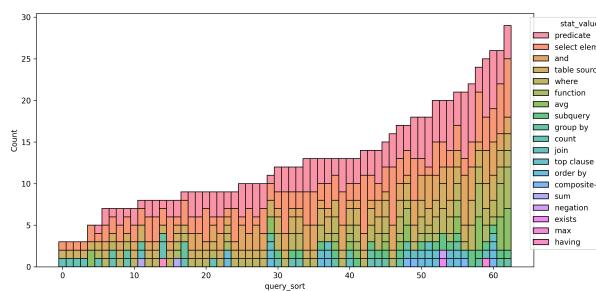


Figure 21: Gold query clause composition - NYSED database

NTSB Database Technical Details.

- Data source format: Microsoft Access
- Migration method: SQL Server Migration Assistant
- Table count: 27
- Column count: 423
- Mean columns per table: 15.67
- NL Questions: 63
- Combined naturalness level: 0.68

A.1.9 Localized Demo Databases Now Available for SAP Business One.

Data Description. The SBOD database [45] is sourced from a publicly available SAP demonstration and training database. It is the largest, and least natural, database within our dataset. Given its schema size, we divided it based on SAP module, and further reduced the schemas used in our benchmark based on the training database cardinality (e.g. we removed most tables containing 0 tuples).

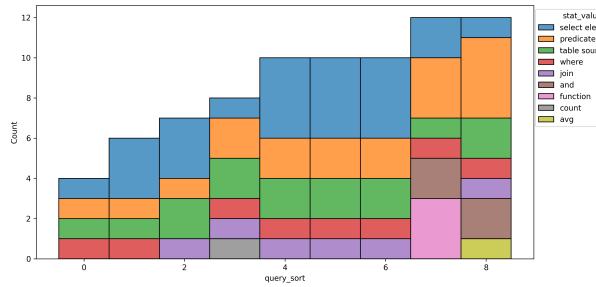


Figure 22: Gold query clause composition - SBOD database module example

SBOD Database Technical Details.

- Data source format: MS SQL Server backup (.bak) file
- Migration method: MS SQL Server backup recovery
- Table count: 2,588
- Column count: 90,477
- Mean columns per table: 34.96
- NL Questions: 100
- Combined naturalness level: 0.49

Module	Tables	Columns	Questions
Banking	40	1720	10
Business Partners	40	1443	10
Finance	61	1988	10
General	71	1035	10
Human Resources	28	440	20
Inventory and Prod.	65	1942	10
Reports	40	734	10
Sales Opportunities	20	283	10
Service	40	875	10

Table 4: SBO Demo Module Schemas

SAP Business One Additional Details. Business One is an enterprise resource planning (ERP) system created by German software systems developer. It is a common platform in government and commercial domains where large-scale business management solutions are required. The SBOD database contains a significant number of tables and columns. Such a schema size poses a problem for generating schema knowledge representations in zero shot prompts, even for large context window variations of evaluated LLMs. To overcome this constraint, we divide the SBOD schema into 9 sub-modules based on schema descriptions published by an online community of SAP practitioners [2]. We further prune the SBOD schemas using the cardinality of the training database, where tables without data entries were excluded from NL questions and prompt schema knowledge.

A.2 NL Questions and Gold Queries

The NL question - SQL query pair artifact consists of 9 .sql files containing between 40 and 100 entries each. Question and query pairs are written in executable .sql files. Natural language questions are written as SQL comments; and SQL is written in the T-SQL dialect employed in MS SQL Server. For public repositories storing the questions, we store them in .zip files in order to reduce the possibility of inclusion in language model training material. Each file is associated with a database in the SNAILS schema collection. Some NL questions contain hint and note entries annotated as HINT and NOTE in lines that follow the NL query. We used neither the hints nor columns in any of the experiments described in this paper, but retain them for possible use in future research.

While we store the data in .sql file format for readability and ease of use, we also offer a NL question loading script (`load_nl_questions.py`) in our repository. This script performs rudimentary parsing of the .sql files and returns a Pandas DataFrame and optional .xlsx formatted spreadsheet.

NL Question - Query Example 1, ASIS Database Question 8. The focus of this benchmark dataset is on the evaluation of schema linking. As such, we were generous with value descriptions, providing literal value strings (e.g. ASIS_HERPS_20H location code in example 1) in the prompt.

```
-- 8: show how many minnows of each stage were counted
      at the location ASIS_HERPS_20H
SELECT stage, sum(count) minnowCountSum
```

```
FROM tblFieldDataMinnowTrapSurveys
WHERE locationID = 'ASIS_HERPS_20H'
GROUP BY stage
;
```

NL Question - Query Example 2, NTSB Database Question 13. Example 2 shows additional code value hints provided in the NL question. In order to enable the recall evaluation statistic, we limited the use of columns and tables in gold queries to the minimum necessary to form a correct query. In the cases where any arbitrary column as an argument in the count function will yield the same result as the *, we use only the * symbol. This eliminates the recall penalty for models selecting an arbitrary column within the count function.

```
-- 13: How many vehicles are there where drugs were present
      (presence code value is 1) and the vehicle was towed
      for a reason not due to disabling damage (towed code is 3)
SELECT COUNT(*) VEHCOUNT
FROM GV
WHERE PARDRUG = 1 AND TOWED = 3
;
```

NL Question - Query Example 3, SBOD Database Human Resources Module Question 18. Questions vary in their complexity. This example shows one of the more complex questions that require multiple projections and joins as well as a selection.

```
-- 18: Show the professional status and educational
      statuses as well as the home and work street
      numbers of employees on the purchasing team.
SELECT StatusOfP, StatusOfE, StreetNoW, StreetNoH
FROM OHEM_employees
JOIN HTM1_teamMembers
  ON employees.empID = teamMembers.empID
JOIN OHTM_emplTeams
  ON teamMembers.teamID = emplTeams.teamID
WHERE emplTeams.name = 'Purchasing'
;
```

A.3 Benchmark Naturalness Comparisons

Though we believe the quality of existing benchmarks is excellent, and the hard work of researchers associated with those projects has resulted in significant improvements in NL-to-SQL system design, we find that a large proportion of these benchmark schemas are canonical, small in composition, and highly natural compared to databases and data sets often encountered in real world scenarios. Using our naturalness classifier described in the main report and the appendix section B, we determine the naturalness levels of the Spider [61] and Bird-SQL [29]. Our classifier model indicates that both Spider and Bird database schemas are highly natural, more so than any real-world schema we acquired for our benchmark (see Figure 23 for a visual comparison). Additionally, we evaluate the naturalness of the real-world schema identifiers in the SchemaPile [11] and find that the SNAILS naturalness better-aligns with SchemaPile than the previously-mentioned benchmarks. Figure 5

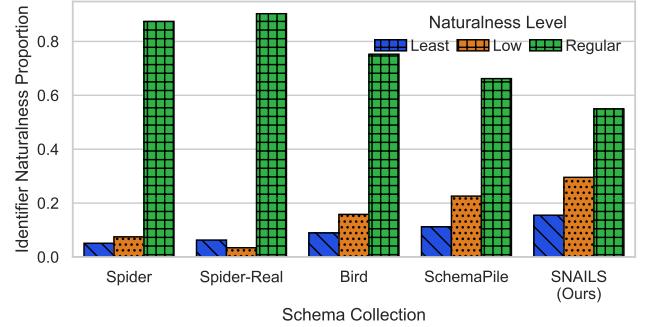


Figure 23: Spider [61] and Bird [29] benchmarks classified with the Davinci-based classifier reveals that both benchmark databases have highly natural identifiers even compared to the most natural of the databases in our proposed benchmark. Our benchmark more closely aligns with the naturalness of the real-world schema collection SchemaPile [11]

B NATURALNESS CLASSIFICATION

B.1 Heuristics-based scoring

Prior to experimenting with ML classifiers, we used a set of heuristics to score the naturalness of each identifier. Comparisons between the heuristics-based scoring approach and ML classification reveals that ML is superior in terms of recall, precision, and F1. We include a description of the heuristics here for completeness, but exclude them from the main body of the report.

- Vectorize an English word vocabulary as frequency counts of letters in the word.
- With a given database identifier, vectorize the identifier as frequency counts of letters in the identifier and downsample to the English word vocabulary to words that have a superset of the letters in the identifier.
- Further downsample the candidate words to words where the letters appear in the same order as the words in the identifier.
- For each word in the downsampled vocabulary, compute the Levenshtein distance between the word and the identifier. This number is called the edit distance.
- For each word, count the number of possible word candidates within 1 and 2 Levenshtein distance from the word. We call this number candidate ambiguity.
- The distribution of candidate ambiguity across our vocabulary is highly skewed, so we take the log of the candidate ambiguity to normalize the distribution.
- We then calculate the naturalness score as the weighted mean of the inverse of the edit distance and the inverse of the log of the candidate ambiguity. This yields values ranging from 0 to 1, where 0 is least natural and 1 is most natural.

B.2 ML Classifier-based scoring

The use of pre-trained language models is a SoTA approach for classification problems [32], and we experiment with various attempts at model-based scoring, including few-shot learning via the

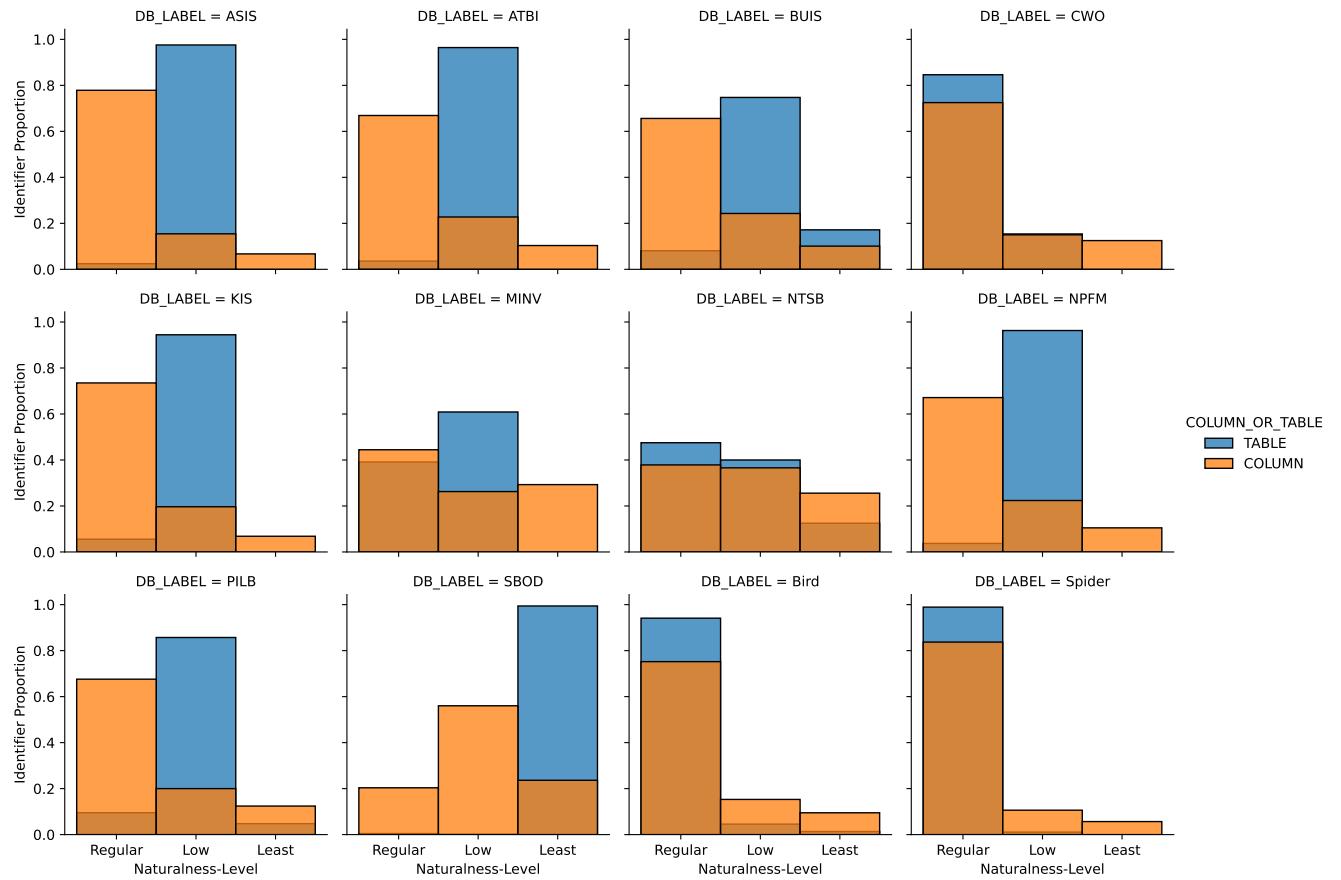


Figure 24: At the individual database schema level, the SNAILS database collection has a diverse arrangement of naturalness levels.

Model	Accuracy	Precision	Recall	F1
GPT-3.5-FewShot	0.646	0.623	0.638	0.630
CANINE-Seq C1	0.719	0.699	0.727	0.712
GPT-4-FewShot	0.742	0.742	0.792	0.766
CANINE-Seq+TG C1	0.829	0.829	0.838	0.833
GPT-3.5-FineTune	0.899	0.878	0.877	0.878
GPT-3.5-FineTune+TG	0.896	0.896	0.897	0.896
CANINE-Seq+TG C2	0.896	0.896	0.898	0.897

Table 5: Performance comparison of different language models for classifying a database identifier's naturalness

GPT API, and finetuning several BERT-like language models on our

dataset of database identifiers to create a second larger collection of identifier naturalness scores. Since the presence of acronyms and abbreviations is a significant determining factor of identifier naturalness, a primary consideration for our naturalness scoring task is the granularity of the tokenizer output. For this reason, we use models that employ either character-level tokenization, word part tokenization, or byte pair tokenization techniques. We select 2 approaches: 1) Use of a foundational LLM in various capacities; and 2) Finetuning of a character-level token language model.

Training data. For finetuning tasks, we train language models using database identifiers extracted from the schemas in the SNAILS real-world database collection. We begin with a human-classified collection (Collection 1); then we employ classifier models trained on Collection 1 to generate a larger set (Collection 2) of machine-classified and human-curated identifier classifications.

Collection 1 The full dataset contains 1,648 manually classified unique schema identifiers. The identifiers are hand labeled as one of 3 naturalness levels (Regular, Low, Least). We randomly divide the data into a training, validation, and test set. This resulted in a distribution of 959 identifiers used for training, 356 for validation, and 333 for testing.

Collection 2 The labeled data set contains 13,722 distinct column identifiers and 3,504 distinct table identifiers for a total size $n = 17,226$. We employ GPT’s Davinci model finetuned on Collection 1 to generate the preliminary naturalness scores. The authors reviewed, and where necessary, modified the model-generated identifier scores to affirm the accuracy of the naturalness classifications. Of the 17,223 identifiers reviewed, 15,527 naturalness scores, or 90.1 percent, were accurately predicted by the Davinci-based model. We manually scored the incorrectly-predicted 1,696 identifiers. For model finetuning, we randomly split the resulting data into training ($n = 10,327$), validation ($n = 3,457$), and test ($n = 3,457$).

B.3 Character Tagging Feature

We include a pre-processing step that generates a sequence of special characters that correspond to the type of characters of the identifier to be classified. The sequence is then concatenated with the identifier and passed to the language models during training and inference. We refer to this approach as *character tagging*, and models employing tagging are labeled with TG in Table 5. Both GPT- and CANINE-based models exhibit improvement in F1 scores using character tagging.

There are intuitive structural differences between abbreviated words and their complete counterparts. Specifically, we observe that word abbreviations generally contain more consonants than vowels, as vowels seem more likely to be removed during abbreviation. We are unsure of a language model’s ability to make use of this observation, so we offer some assistance in the form of a pre-processing step that generates a sequence of special characters that correspond to the type of characters of the identifier to be classified. We concatenate the tag sequence to the identifier and pass it to the language models for training and inference. Special characters include:

- ^: Vowels
- +: Consonants
- #: Numbers
- \$: Special characters
- *: Any character not in the above categories

We refer to this approach as *character tagging*, and models employing tagging are labeled with TG in Table 5 of the main report.

For example, the identifier AuthorID_5 would be pre-processed as follows:

AuthorID_5 ^++^+++\$#

Both GPT- and CANINE-based models exhibit improvement in F1 scores when using the character tagging feature.

B.4 GPT 3.5/4 Turbo Few-Shot-Based Scoring

We experiment with the effectiveness of few-shot prompting to classify identifier naturalness. We opt to provide one set of instructions followed by a series of 25 randomly selected human-validated examples of naturalness levels. We perform text replacement on the trailing row by replacing the _IDENTIFIER_ text with the identifier to be classified. This approach does not require model pre-training; but this convenience is paid for in terms of the number of tokens in the prompt, and classifying a large schema with this method can incur rather high LLM usage costs.

The following is a list of database identifiers and labels that indicate how closely they resemble natural english words:

N1: most natural english words
 N2: second most natural english words
 (e.g. abbreviations or combinations of natural words and acronyms)
 N3: third most natural english words
 (e.g. very short abbreviations with obscured meaning or acronyms)

```
identifier: CASENO Label: N1
identifier: BENTHOS_TotalAreaSampled_m2 Label: N2
identifier: CAUSE3 Label: N1
identifier: MT_RIVPACS_2011_OTU Label: N3
identifier: ACTIVATE Label: N1
identifier: MotorcycleChassisTypeId Label: N1
identifier: First_Name Label: N1
identifier: IPCAREA_2ND Label: N2
identifier: INJNO Label: N2
identifier: tbl_MicroHabitat Label: N2
identifier: EMMSGCSEYE Label: N3
identifier: HEADRESTDAM Label: N2
identifier: AutoPedestrianAlertingSound Label: N1
identifier: ModelTest Label: N1
identifier: tlu_topo_position Label: N2
identifier: Understory_Comp Label: N1
identifier: BAGDAMAGE Label: N1
identifier: HARNESSDESIGN Label: N1
identifier: Coord_Syst Label: N2
identifier: CINJSEV Label: N2
identifier: JKWTG12 Label: N3
identifier: _IDENTIFIER_ Label:
```

B.5 GPT Davinci Fine Tuning

We train the Davinci-based completion models using the OpenAI command line API. We generated models with character tagging, as well as models without tagging. Below is an excerpt from the tagging-based training data.

```
{"prompt": "ADDRESS ^++^++ ->", "completion": " N1"}  

>{"prompt": "AIS ^+ ->", "completion": " N3"}  

>{"prompt": "AISCODE ^++^+^ ->", "completion": " N3"}  

>{"prompt": "BACKBPILL +++++^++ ->", "completion": " N2"}  

>{"prompt": "ALIGNMENT ^++^++^++ ->", "completion": " N1"}  

>{"prompt": "ARRMEDICAL ^++^+^+ ->", "completion": " N2"}
```

Inference using character tagged models requires appending the tag to the identifier in the same format as the training data.

B.6 CANINE Fine Tuning

CANINE [8] is a BERT-based language model that tokenizes inputs at the character level. We trained a sequence classification head using both generation 1 and generation 2 data sets using a single NVIDIA GTX 1080 GPU. We employed the HuggingFace Transformers library, CUDA 12.1, and Torch 2.0.1 to fine tune the ‘google/canine-s’ model. Hyperparameter tuning was conducted using the optuna library, which resulted in the parameter settings:

- Optimizer: adamw hf
- Learning rate: 4.910828967396573e-05
- Per device training batch size: 24
- Per device evaluation batch size: 12
- Number of epochs: 15
- Weight decay: 0.04168784348465411

Models were trained both with, and without, the character tagging feature. We offer the evaluated models in our project repository. The *ml-naturalness-canine-sequence-classifier-inference.py* Python file contains the *CanineIdentifierClassifier* class. This class provides a simple *classify_identifier* method for using the CANINE model to classify identifier naturalness with or without character tagging.

C NATURALNESS-MODIFIED IDENTIFIERS

C.1 Dataset Naturalness Classifications

Identifier naturalness within each dataset is categorized using the *N1* (Regular), *N2* (Low), and *N3* (Least) categories. Naturalness of table and column identifiers are cataloged both separately, and in consolidated form (i.e. tables and columns together). Additionally, we calculate a combined naturalness score for consolidated identifiers using category weights.

$$\text{CombinedNaturalness} = 1.0 * \text{Regular} + 0.5 * \text{Low} + 0.0 * \text{Least} \quad (5)$$

where *N1*, *N2*, and *N3* are proportions of schema identifiers in each respective category within the total count of identifiers in the identifier's source schema.

C.2 Training Data Collections

For fine tuning tasks, we train language models using database identifiers extracted from the schemas in the SNAILS real world database collection. We begin with a human-classified collection (Collection 1); and employ classifier models trained on Collection 1 to generate a larger set (Collection 2) of machine-classified and human-modified identifier classifications.

Collection 1 The full data set contains 1,648 manually classified unique schema identifiers. The identifiers are classified as one of three naturalness levels (Regular, Low, Least). We randomly divide the data into a training, validation, and test set. This resulted in a distribution of 959 identifiers used for training, 356 for validation, and 333 for testing.

Collection 2 The training data set contains 13,722 distinct column identifiers and 3,504 distinct table identifiers for a total size $n = 17,226$. We employ a fine tuned model using GPT-3.5 Davinci to generate the preliminary naturalness scores. The authors reviewed, and where necessary, modified the model-generated identifier scores to affirm the quality of the training, validation, and test data. Of the 17,223 identifiers reviewed, 15,527 naturalness scores, or 90.1 percent, were accurately predicted by the Davinci-based model. We manually re-score the incorrectly-predicted 1,696 identifiers. For model fine tuning, we randomly split the resulting data into training ($n = 10,327$), validation ($n = 3,457$), and test ($n = 3,457$).

Naturalness modification is the process of changing an identifier in such a way that it assumes a naturalness category that is not

its original classification (see Figure 25). Modifying an identifier to become less natural is useful for creating benchmark schemas of varying naturalness levels. The same benefit applies to the process of modifying less natural identifiers to become more natural; this direction of modification also generally yields improved NL-to-SQL performance, as is demonstrated in the experiment and evaluation sections of this report.

Naturalness-modified identifiers generated by the ML-based approaches described next are human-validated and, when necessary, modified. Once validated, the identifiers are added to our ground truth dataset and used for the prompt and query naturalness modification processes described elsewhere in this report.

C.3 Decreasing Naturalness (Abbreviation)

Decreasing naturalness generally involves the removal of characters from an identifier in a manner that shortens the length while retaining some structure that still allows for some measure of readability. This cannot be achieved by randomly removing characters from an identifier; so we elect to use machine learning-based approaches to decrease identifier naturalness. As with our classification approaches, we experiment with both engineered few shot prompts targeted at a general purpose foundational LLM (GPT), and a fine tuning approach (GPT Davinci).

PPT Davinci Fine Tuning Abbreviation. Separate models are trained for conversion tasks from one naturalness level to a model-specific alternative naturalness level. This resulted in the following fine tuned models:

- Regular to Low
- Regular to Least
- Low to Least

Each fine tune dataset consists of 176 randomly selected identifiers and human-created naturalness modifications. Below is an example of Regular to Least model training data:

```
{"prompt": "Plot ->", "completion": " p\\n"}  
{"prompt": "Metals ->", "completion": " mt\\n"}  
{"prompt": "Station_ID ->", "completion": " S_ID\\n"}  
{"prompt": "FUELEAK ->", "completion": " F_Lk\\n"}
```

The outputs of these finetuned models require significant adjustment by human researchers; so we elect to employ an alternative approach described next.

GPT Few Shot Abbreviation. GPT 3.5-based few shot prompting (see example below) resulted in the most consistent outputs, with a reasonably low prompt token count. Rather than explaining the different categories followed by an instruction to convert an identifier to a specific category, we find that providing a simple instruction to *abbreviate the database schema identifier to make it slightly shorter* followed by several examples is more effective.

Abbreviate the database schema identifier to make it slightly shorter:

Protocol_Name -> Protcl_Nm

Abbreviate the database schema identifier to make it slightly shorter:
WaterTemperature -> WaterTemp

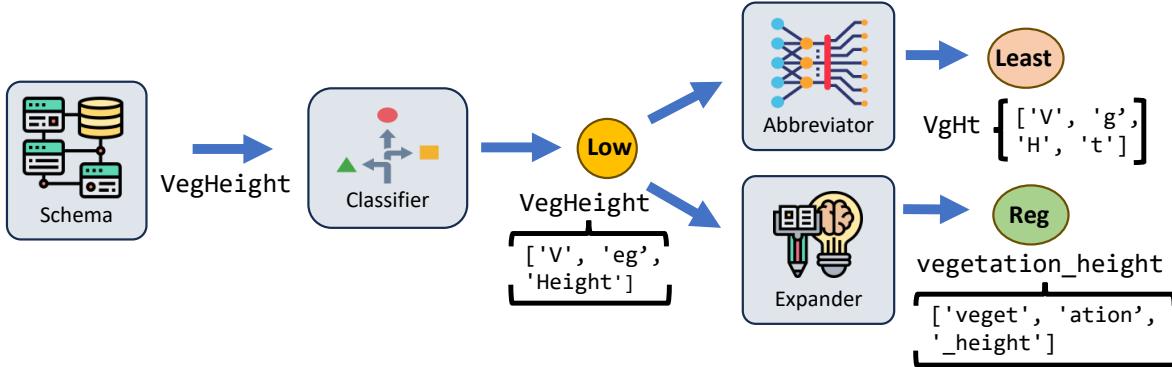


Figure 25: Schema identifiers in our benchmark dataset are classified into a naturalness category and modified as appropriate. Modified identifiers comprise the schema crosswalks used for schema modification during experiment query inference.

Abbreviate the database schema identifier to make it slightly shorter:
Customer -> Custmr

Abbreviate the database schema identifier to make it slightly shorter:
IDENTIFIER ->

C.4 Increasing Naturalness (Expansion)

Increasing naturalness requires the expansion of an abbreviated identifier. A recent attempt at performing this task with model fine tuning [63] was made during our research; and it appears to be a promising direction for research. However, we elect to enrich our process with external database metadata.

Expansion Process. In order to accomplish this, we engineered a solution that employs a database metadata reader capable of reading .pdf, .xml, and .csv formatted metadata. Metadata is read and indexed at the word level, where an array of file locations (page and line numbers for pdf, line numbers for xml and csv) where words occur are mapped to each word. When a user keys in an identifier to modify, file locations where the identifier exists in the document are returned via index lookup. These index locations are used as the centerpoints of context windows that retrieve the surrounding content. This content is added to a fewshot prompt to provide the language model with document content that is likely to contain references to, and definitions of, the provided identifier.

The fewshot prompt for generating an expanded identifier adheres to the template:

Using the following text extracted from a data dictionary:
__CONTEXT__

In the response, provide only the old identifier and new identifier (e.g. "old_identifier, new_identifier").

Create a meaningful and concise database identifier using SQL compatible complete words to represent abbreviations and acronyms for only the identifier __IDENTIFIER__:

The __CONTEXT__ placeholder is replaced with up to ten context window-length excerpts from the database metadata. This is an example of a completed prompt using the NYSED .pdf based data manual with context window of 200 characters:

Using the following text extracted from a data dictionary:
r school Text 255
YEAR Reporting Year (2021 for 2020 -21; 2022 for 2021 -22)
Number 4
NUM_TEACH Number of teachers as reported in the Student Information Repository System (SIRS) Number 12
NUM_TEACH_INEXP Number of teachers with fewer than four years of experience in their positions Number 12
PER_TEACH_INEXP Percent of teachers with fewer than four years of experience in their posi

In the response, provide only the old identifier and new identifier (e.g. "old_identifier, new_identifier"). Create a meaningful and concise database identifier using SQL compatible complete words to represent abbreviations and acronyms for only the identifier num Teach_inexp:

num Teach_inexp, number_of_teachers_inexperienced

In this successful example, we see that the identifier num Teach_inexp has been expanded to a more natural number_of_teachers_inexperienced. This is despite the observation that the data retrieved from the .pdf file is quite unstructured and contains document artifacts. A sufficiently wide context window coupled with the retrieval of multiple occurrences of the identifier in the document generally results in valid expansions.

Prompt Building. In order to generate few shot prompts over an arbitrary metadata source, some prompt engineering is necessary. Generally, hand-crafted prompt building is suitable approach; but it does not scale nor does it lend itself to an automated solution that can be deployed beyond a research lab. To make this process more portable, we introduce a command line-based subroutine that enables the automatic build of a five example few shot prompt. In this process:

- (1) User enters an identifier
- (2) Zero shot prompt -> expanded identifier
- (3) User reviews and validates identifier
- (4) Correct: identifier added to example list
- (5) Incorrect: User tries again with different identifier
- (6) User enters another identifier
- (7) Few shot prompt (with prior successes as examples) -> expanded identifier
- (8) Correct: identifier added to example list
- (9) Incorrect: User tries again with different identifier
- (10) Process repeats until five successful examples are generated

Once a fewshot prompt has been created for a given database’s metadata, the prompt is stored for any future program runs. This particular aspect of our project was built to support our research efforts; and we did not perform any experiments to evaluate its overall accuracy and usability. We leave these tasks as future research opportunities.

D NL-TO-SQL BENCHMARKING SETUP

D.1 Prompting

Prompts are generated dynamically during inference runtime and include schema knowledge, task and syntax instructions, and a natural language question. Schema information is extracted from the target database system tables and encoded into relational diagram format as schema knowledge. Instructions include SQL dialect type, and answer format (e.g. provide only a SQL query without explanation).

D.2 Prompt Naturalness Modification

Prompt naturalness modification is necessary when generating SQL queries over schemas with modified identifiers. In order to prevent producing additional database instances with renamed identifiers, we employ a middleware approach where modified identifiers are retrieved from a mapping of native identifiers to the target naturalness level. The prompt is generated using native schema identifiers, and table and column names are encased in XML-like opening and closing `<TABLE_NAME>` and `<COLUMN_NAME>` tags. Native identifiers and their tags are replaced with modified identifiers using standard Python string replacement (e.g. `str.replace(target, value)`).

Prompt naturalness conversion example. The objective is to modify the naturalness of the Klamath Invasive Species (KIS) schema to the least natural form. The first step is to generate a tagged prompt which is formed using the database metadata accessed via system tables. A tagged prompt table with columns and datatypes takes on the form:

```
#<TABLE_NAME>tlu_Species_WHIS</TABLE_NAME>
```

```
( <COLUMN_NAME>Species</COLUMN_NAME> nvarchar,
<COLUMN_NAME>SampleYear</COLUMN_NAME> nvarchar,
<COLUMN_NAME>Park</COLUMN_NAME> nvarchar
)
```

Each table and its columns occupies a single line within the prompt. The resulting prompt after string replacement appears as the following in the final prompt:

```
#TSW( Sp nvarchar, S_Yr nvarchar, Pk nvarchar)
```

The modified schema knowledge is presented to the LLM as though it is the native schema.

D.3 NL-to-SQL Prediction

Four LLMs are integrated into the experiment pipeline: GPT 3.5 Turbo, GPT 4, Gemini 1.5 Pro, and Phind Code Llama. In earlier experiments, we also used CodeLlama 34b, CodeLlama 7b, and Google Bison. Due to these models being superseded by more-capable variants, we exclude them from the results in our main report.

GPT 3.5 Turbo and GPT 4. GPT-based [36] generations use the `gpt-3.5-turbo-16k` model accessed using OpenAI’s API services [37]. GPT-based models employ BPE using tiktoken [38] to tokenize inputs and decode model outputs.

Due to the size of larger schema knowledge prompts, we make use of the GPT 3.5 model [36] with a context window of 16,000 tokens `gpt-3.5-turbo-16k`. For consistency, all queries were generated with a 0 temperature, 1 top p, 0 frequency penalty, 0 presence penalty. Responses are fetched from the OpenAI Python `ChatCompletion` class’ `create` method with a single message (prompt) passed in the `user` role.

Code-Bison. Code-Bison is based on Google’s PaLM 2 LLM. Although model architecture details are not readily available, we deduce from the PaLM 2 tech report [5] that the model uses the BERT tokenizer which relies on WordPiece tokenization [10] to encode inputs.

Queries were generated with a 0 temperature, and 1 top p. The largest available context window was the experimental 32k model `code-bison-32k`, which we use due to schema knowledge token counts. Generated queries are prepended by a ““sql”” label, which we remove with standard Python string replacement.

Gemini 1.5. Gemini is Google’s most recent foundational model which exhibits superior performance across a wide range of tasks. We derive from the Gemma technical report [55] (an open source LLM based on the same architecture as Gemini) that it uses a SentencePiece tokenizer. Of note, the context window of the Gemini model is one million tokens, which allows the inclusion of the entire schema of the extremely large SBODemoUS schema.

Code-Llama. Two variants of Meta’s open source Code Llama [?] base model are used: CodeLlama-7b and CodeLlama-34b. The 7b model was deployed on an AWS EC2 p3.8xlarge instance equipped with 4 NVIDIA Tesla V100 GPUs, each with 16GiB GPU memory (64GiB total). The 34b model was deployed on an AWS EC2 p3.16xlarge instance equipped with 8 Tesla V100 GPUs (128GiB total). Both deployments utilize the VLLM [26] large language model

deployment infrastructure to host an API endpoint integrated into the query generation and evaluation pipeline. Llama model tokenizers employ BPE for input and output encoding.

Queries were generated with a 0 temperature, 1 top p, 0 frequency penalty, and 0 presence penalty. By default, Code Llama provides a size 100,000 token context window. The Code Llama model used in our experiment is a completion model, so it was necessary to append the SELECT keyword to the end of the prompt prior to inference.

Phind Code-Llama. Phind Code-Llama [44] is a fine-tuned version of the 34b parameter Code Llama model. We used the TogetherAI API to access this model. We find that the finetuning appears to have improved its performance as compared to the baseline Code Llama 34b version, and as such we replace our prior analysis using Code Llama 34b with the results generated using the Phind model.

D.4 Query Naturalness Modification

When a query is formed against a schema with naturalness-modified identifiers, it is necessary to replace the modified identifiers with the native identifiers prior to executing the query over the target database. Simple string replacement is not sufficient in this case, because some identifier names may sometimes be substrings of other identifiers within a query; and replacing one may corrupt the other. We employ a Java-based parser and AST generator [3, 43] to build a parser system for tagging table and column identifiers in a query.

Query naturalness modification example. To answer the question *For each location type, show a count of locations in shasta county*, GPT 3.5 generated the query with lowest naturalness schema knowledge:

```
SELECT LcTp, COUNT(*) AS LocationCount
FROM Locs
WHERE Cty = 'Shasta County'
GROUP BY LcTp
```

This query is converted to all capital characters and passed to the query parser tagger via API. The tagger traverses the AST using a listener class to encase tables and columns with identifier opening and closing tags. The parser also returns a list of aliases generated within the query. This allows consuming systems to ignore tagged aliases within the tagged query:

```
@BEGINAGGREGATE
SELECT
    <COLUMN_NAME> LCTP </COLUMN_NAME> ,
    COUNT ( * ) AS LOCATIONCOUNT
FROM <TABLE_NAME> LOCS </TABLE_NAME>
WHERE
    <COLUMN_NAME> CTY </COLUMN_NAME>
    =
    'SHASTA COUNTY'
GROUP BY
    <COLUMN_NAME> LCTP </COLUMN_NAME>
@ENDAGGREGATE
@BEGINALIASES
<COLUMN_ALIASES>LOCATIONCOUNT</COLUMN_ALIASES>
@ENDALIASES
```

With the tagged query and aliases, the query is converted into a form that is compatible with the target native schema. This is

accomplished using the schema identifier mapping dataset and standard Python string operations (e.g. str.replace("<TABLE_NAME> LOCS </TABLE_NAME>", "TBL_LOCATIONS")).

@DENATURALIZED RESPONSE:

```
SELECT
    [LOC_TYPE],
    COUNT (*) AS LOCATIONCOUNT
FROM [TBL_LOCATIONS]
WHERE [COUNTY] = 'SHASTA COUNTY'
GROUP BY [LOC_TYPE]
```

This completes the query modification step. The modified query is then used to extract results from the target database for result set matching evaluation.

E PERFORMANCE EVALUATION

E.1 Query Execution

Predicted and gold queries are executed over target databases with native schemas using the PyOdbc Python library. Valid result sets are stored as Pandas Dataframes for comparison. Invalid queries generate MS SQL Server error messages; messages are captured for categorization.

E.2 Execution Result Set-Subset Matching

Result set and subset comparison consisted of a comparison of the columns and rows returned by the two queries. Because it is possible for semantically equivalent query results to differ in terms of column ordering, column and aggregate function aliasing, row ordering (when an order is not specified in the question), and even the number of columns returned (as long as the predicted column set is a superset of the gold column set), result set comparison was performed using a series of rules implemented in Python.

Result cardinality The number of tuples, denoted as $|C_G|$ and $|C_P|$ in the predicted result R_P and the gold result R_G must be equal, and must be greater than 0:

$$\forall C_G \in R_G, C_P \in R_P (|C_G| = |C_P|) \wedge (|C_G| > 0 \wedge |C_P| > 0)$$

Empty sets are tagged as undetermined and retained for syntactic comparison. Non-empty, non-equal-size sets are tagged as non-matches and withheld from further analysis.

Projection completeness The columns in the predicted result R_P must be a superset of the columns in the gold result R_G . Columns C_P and C_G equivalence is determined by value comparisons between tuples T_G and T_P of the columns' contents:

$$\forall C_G \in R_G, \exists C_P \in R_P \text{ such that } \forall T_G \in C_G, T_G \in C_P$$

Column match candidates are determined via pairwise comparison of the sorted values in each column in R_G to each column in R_P . Full result sets R_G and R_P are sorted by corresponding column match candidates, with columns containing the most unique values serving as the primary sort key. With both R_G and R_P sorted by column match candidates, the two sets are compared row-wise for all columns in the set of column match candidates. If the two sets are not equal, the result sets are considered semantically non-equivalent.



Figure 26: Screenshot of the query manual validation tool. This example depicts a query that passed set-subset matching, and is currently classified as ungraded. Helper information indicates that although the results matched, the incorrect table was selected during inference (AHEM instead of OHEM). This example was classified as incorrect.

E.3 Human Evaluation

Predicted queries that pass execution result set-subset matching are further evaluated by a human researcher to rule out false positives. Predictions that fail result set comparison are pre-classified as also failing manual matching. Predictions that pass set comparison are classified as ungraded until reviewed by a human researcher. Once reviewed using the Python-based GUI evaluation tool (see Figure 26), the predicted query receives its final matching score. Human evaluation resulted in scoring as incorrect 41 predictions that passed result set matching which is approximately two percent of all result set matched queries.

E.4 Schema Linking End-to-end Example

To better pinpoint schema linking performance, we devise a new approach for evaluating NL-to-SQL generation. In this approach, performance is measured using set comparisons between sets of identifiers within gold and predicted SQL queries. Recall is the primary metric. Precision and F1 are available, but less helpful, due to penalization for additional predicted columns that do not cause a query to be incorrect. The formulae for deriving linking metrics are printed in subsection 5.2 of the main report.

Linking Evaluation Example (ATBI Question 30, CodeLlama-34b). This is an example of a predicted query that fails result set match comparison.

Natural Language Question. Which tree species were recorded as mature overstory but not as saplings? Include the species name and common name.

Gold Query. The gold query is quite complex, where using EXISTS clauses is a possible approach for solving.

```
SELECT species, CommonName
FROM tlu_PlantSpecies sp
WHERE EXISTS(
```

```
SELECT overstory_id
FROM tbl_Overstory
WHERE SpCode = sp.SpeciesCode
) AND NOT EXISTS (
SELECT Seedlings_ID
FROM tbl_Seedlings
WHERE SpCode = sp.SpeciesCod
)
```

Predicted Query (Code Llama 34b). As we can see below, Code Llama 34b did not provide a correct solution. However, failure of result set matching does not preclude us from gleaning useful information from the result.

```
SELECT DISTINCT
tlu_PlantSpecies.genus,
tlu_PlantSpecies.subgenus,
tlu_PlantSpecies.species,
tlu_PlantSpecies.subspecies,
tlu_PlantSpecies.SpeciesCode,
tlu_PlantSpecies.CommonName
FROM
tlu_PlantSpecies
LEFT JOIN tbl_Overstory
ON tbl_Overstory.SpCode
= tlu_PlantSpecies.SpeciesCode
LEFT JOIN tbl_Saplings
ON tbl_Saplings.SpCode
= tlu_PlantSpecies.SpeciesCode
WHERE
tbl_Overstory.SpCode IS NOT NULL
AND tbl_Saplings.SpCode IS NULL
ORDER BY
tlu_PlantSpecies.genus,
tlu_PlantSpecies.subgenus,
tlu_PlantSpecies.species,
tlu_PlantSpecies.subspecies,
tlu_PlantSpecies.SpeciesCode,
tlu_PlantSpecies.CommonName
```

Identifier Set Extraction. Using the purpose-built T-SQL parser, we extract a set of identifiers from the gold and predicted queries. Note that we do set comparison; so although columns are referenced in multiple clauses in the predicted query, we only measure the presence of a column or table once.

With a set QI_g of identifiers present in the gold query and a set of identifiers QI_p present in the generated (or predicted) query, we calculate recall, as well as f1 and precision.

Gold query identifiers QI_g :=

```
{
'TLU_PLANTSPECIES', 'TBL_OVERSTORY', 'TBL_SEEDLINGS',
'SPECIES', 'SPECIESCODE', 'COMMONNAME', 'SPCODE',
'OVERSTORY_ID', 'SEEDLINGS_ID'
```

Predicted query identifiers QI_p :=

```
{
'TLU_PLANTSPECIES', 'TBL_OVERSTORY', 'TBL_SAPLINGS',
'SPECIES', 'SPECIESCODE', 'COMMONNAME', 'SPCODE',
```

```
'GENUS', 'SUBSPECIES', 'SUBGENUS'
}
```

Identifier Set Comparisons. True positives are the intersection $QI_g \cap QI_p$ =

```
{
'TLU_PLANTSPECIES', 'TBL_OVERSTORY',
'SPECIES', 'SPECIESCODE', 'COMMONNAME', 'SPCODE'
}
```

$$QueryRecall = \frac{|QI_g \cap QI_p|}{|QI_g|} = \frac{6}{9} = 0.667$$

$$QueryPrecision = \frac{|QI_g \cap QI_p|}{|QI_p|} = \frac{6}{10} = 0.60$$

$$QueryF1 = \frac{2(Recall * Precision)}{Recall + Precision} = 0.632$$

So we see that although the predicted query failed in terms of execution result set comparison, we can still grade it in terms of linking performance. In other words, we can assign *partial credit* to predicted queries where correct schema identifiers are recalled.

E.5 Error Classifications

Microsoft SQL Server error codes were used to classify server-side failure states resulting from the execution of generated SQL queries. Execution errors generally fall into one of two categories: syntax errors, and schema errors. Syntax errors include occurrences such as incorrect SQL syntax variants (e.g. LIMIT instead of TOP), incomplete or absent GROUP BY clauses, and others. Schema errors generally result from language model hallucination of non-existent tables or columns, the use of a column that is not a member of a table referenced in the query, or operations incompatible with a column's data type.

F RESULTS

F.1 Possible Evidence of Schema Exposure

When evaluating individual NL-to-SQL translations, we observe that for the NPFM and KIS schemas, Code-Bison appears to have prior exposure to information about lookup values. In the query below, it correctly forms the predicate “WHERE CanPos = 2”; whereas all other LLMs made incorrect guesses of the same value.

NPFM 15: How many overstory's have a codominant canopy position?

Gold Query (Correct):

```
select count(*) OverstoryCount
from tlu_Can_Pos c
join tbl_Overstory o on c.CanPos_Num = o.CanPos
where CanPos_Name = 'codominant'
;
```

GPT 3.5 (Incorrect):

```
SELECT COUNT(*)
FROM tbl_Overstory
WHERE CanPos = 'Codominant'
```

Code-Llama-7b (Incorrect)

```
SELECT COUNT(Overstory_ID)
FROM tbl_Overstory
WHERE CanPos = 1
```

Code-Llama-34b (Incorrect)

```
SELECT COUNT(*)
FROM tbl_Overstory
WHERE CanPos = 1
```

Code-Bison (Correct)

```
SELECT COUNT(*)
FROM tbl_Overstory
WHERE CanPos = 2;
```

None of the NL questions against the KIS schema were phrased in a way that would reveal examples of prior schema knowledge, so we perform post-hoc schema probing using Code-Bison and GPT 3.5. In this example, we ask the LLM to retrieve enumeration codes from the KIS enumeration lookup table. The question is phrased in a way that requires the use of a predicate applied to the Enum_Group column with a lookup value not present in the prompt. We attempt to avoid coincidental lexical matching by phrasing the group description in a way that would not be reduced to the correct lookup code without prior exposure. As we see below, Code-Bison correctly recalls the value whereas GPT-3.5 makes an incorrect guess.

Show me the codes for the enumeration group that contains the height of vegetation measured in meters.

Code-Bison (Correct)

```
SELECT Enum_Code
FROM tlu_Enumerations
WHERE Enum_Group = 'VegHeight';
```

GPT 3.5 (Incorrect):

```
SELECT Enum_Code
FROM tlu_Enumerations
WHERE Enum_Group = 'Vegetation Height (meters)'
```

F.2 Execution Accuracy

In this section, we provide a more fine-grained analysis of execution accuracy by examining NL-to-SQL accuracy within databases by schema naturalness level. Figure 28 provides a detailed view of accuracy performance while reminding the reader of the naturalness scores of each Native schema. With one exception (Llama-34b over NTSB), databases with native schema scores less than 0.69 exhibit improvement with schemas modified to Regular level naturalness. We also note that for databases with higher naturalness scores, Native schemas generally perform the best.

Figure 27 is a visual representation of the tabular data in Figure 28. This view of the results makes the performance difference between databases and LLMs quite apparent. These between-database execution accuracy differences suggest that other schema-related factors may be at play. Potential challenges could include the use of composite keys (NTSB), very large schemas (SBOD), and the presence of whitespace in native schema identifiers (NYSED).

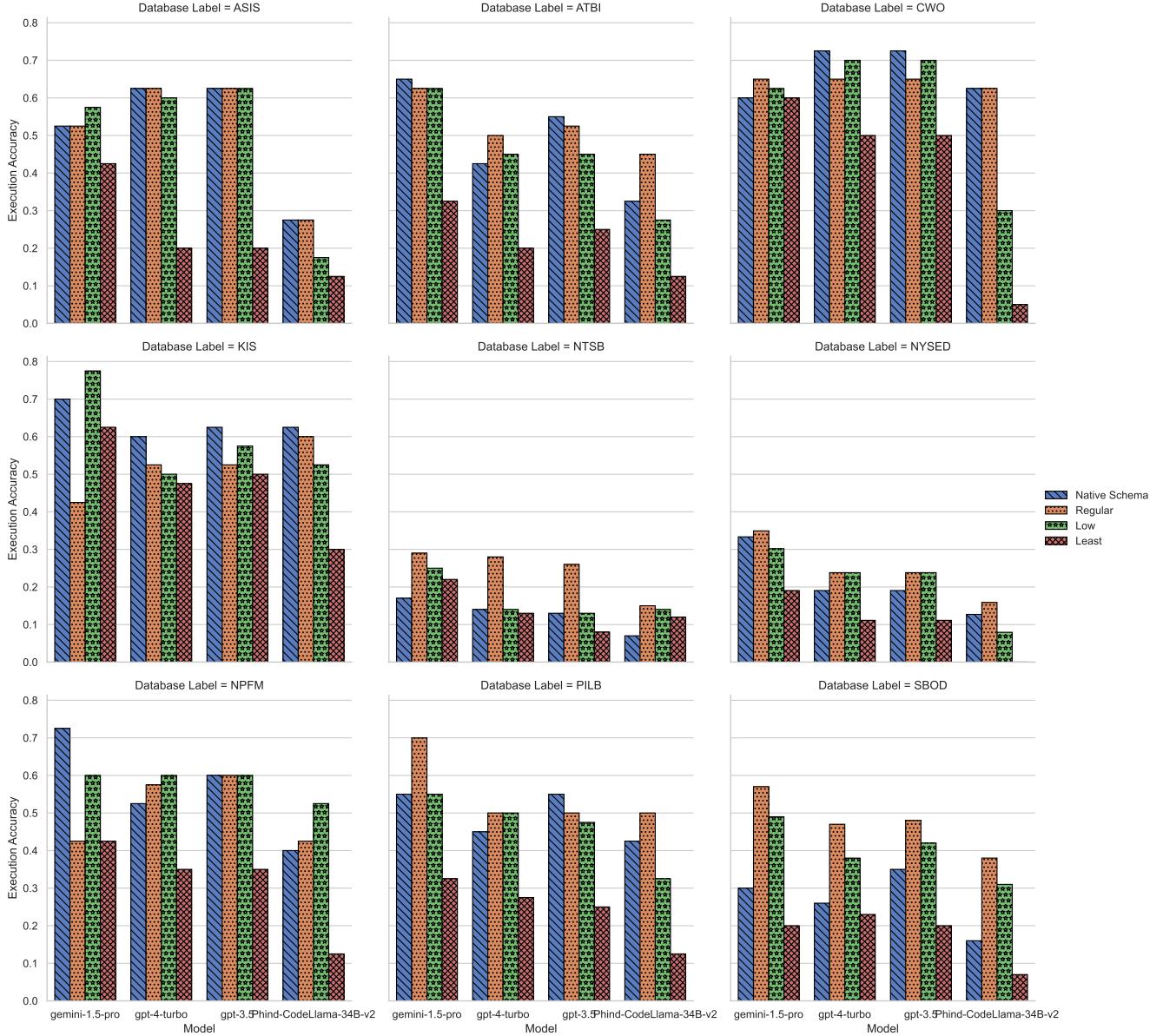


Figure 27: Execution accuracy by database and language model. This view of the results makes the performance difference between databases and LLMs quite apparent. These between-database execution accuracy differences suggest that other schema-related factors may be at play. Potential challenges could include the use of composite keys (NTSB), very large schemas (SBOD), and the presence of whitespace in native schema identifiers (NYSED).

F.3 Schema Linking

Naturalness effect on schema linking. Figures ??, ??, ??, and ?? indicate non-parametric distributions of combined naturalness and linking scores. Additionally, there is not a strong visual indication of correlation; but a linear regression plot suggests some relationship between combined naturalness and schema linking performance.

A Kendall-Tau distribution study (Tables 35a–49b) of correlation between naturalness and schema linking suggests a moderate

and statistically significant relationship between the two for most combinations of language model, and schema type (modified or unmodified). Generally, we see that more natural query identifiers result in higher schema linking outcomes.

Tables 38a–49b indicate that proportions of a naturalness category within a set of query identifiers also have an effect on linking performance. Specifically, as the proportion of Regular naturalness identifiers increases, so does schema linking. We also observe that

Model	Database (Score)	ASIS (0.77)	ATBI (0.70)	CWO (0.84)	KIS (0.79)	NPFM (0.70)	NTSB (0.59)	NYSED (0.68)	PILB (0.75)	SBOD (0.49)
	Category									
gemini-1.5-pro	Native	0.53	0.65	0.60	0.70	0.72	0.17	0.33	0.55	0.30
	Regular	0.53	0.62	0.65	0.42	0.42	0.29	0.35	0.70	0.57
	Low	0.57	0.62	0.62	0.78	0.60	0.25	0.30	0.55	0.49
	Least	0.42	0.33	0.60	0.62	0.42	0.22	0.19	0.33	0.20
gpt-4-turbo	Native	0.62	0.42	0.72	0.60	0.53	0.14	0.19	0.45	0.26
	Regular	0.62	0.50	0.65	0.53	0.57	0.28	0.24	0.50	0.47
	Low	0.60	0.45	0.70	0.50	0.60	0.14	0.24	0.50	0.38
	Least	0.20	0.20	0.50	0.47	0.35	0.13	0.11	0.28	0.23
gpt-3.5	Native	0.62	0.55	0.72	0.62	0.60	0.13	0.19	0.55	0.35
	Regular	0.62	0.53	0.65	0.53	0.60	0.26	0.24	0.50	0.48
	Low	0.62	0.45	0.70	0.57	0.60	0.13	0.24	0.47	0.42
	Least	0.20	0.25	0.50	0.50	0.35	0.08	0.11	0.25	0.20
Phind-CodeLlama	Native	0.28	0.33	0.62	0.62	0.40	0.07	0.13	0.42	0.16
	Regular	0.28	0.45	0.62	0.60	0.42	0.15	0.16	0.50	0.38
	Low	0.17	0.28	0.30	0.53	0.53	0.14	0.08	0.33	0.31
	Least	0.12	0.12	0.05	0.30	0.12	0.12	0.00	0.12	0.07

Figure 28: Execution accuracy by database and language model. Databases with native schema scores less than 0.69 exhibit improvement with schemas modified to Regular level naturalness. We also note that for databases with higher naturalness scores, Native schemas generally perform the best.

in most cases, as the proportion of Low increases, schema linking generally improves, but not to the same degree as for Regular identifiers. The most striking effect comes from the proportion of Least identifiers, where as the proportion of Least naturalness increases, schema linking performance decreases.

Mean token to character ratio effect on schema linking.

Figure 32 suggests that in most cases, there is a moderate correlation between the mean token-character ratio of identifiers in a query, and the recall score. This visual correlation appears for GPT and Llama tokens, but is less apparent for Bison tokens. The Kendall-Tau correlation corroborates this visual indicator (Tables reftable:tokenratio-recall-ktau-native and reftable:tokenratio-recall-ktau-all) which shows us that there is significant evidence that for GPT and Code Llama tokenizers over both native and modified schemas, that as the token-to-character ratio increases, schema linking performance decreases. This same observation does not hold for the Code Bison tokenizer over the native schemas. Because of the inconsistent power of the token-character ratio measurement, we believe it may not serve as a viable proxy for naturalness in all cases.

G RELATED WORK

Ontology Mapping. Schema modifications and intermediate representations to enhance performance in a specific context extend beyond NL-to-SQL applications. Mapping relational database schemas to ontologies is an approach used to improve schema-to-schema integration and web application application-database interfaces [60]. This improves the semantic description of underlying data, which is often a desirable feature in web applications that interact within the semantic web [21]. While ontological mapping

of a relational database can improve performance in this context; we see less evidence that such an approach is useful or necessary in NL-to-SQL applications, though this may serve as a compelling opportunity for future research.

H ADDITIONAL TABLES AND FIGURES

The remaining pages contain several figures and tables of fine-grained analysis of dataset distributions and performance correlations.

Schema Linking Performance Scatter Plots. Figures 29–33 are results of various data explorations and experiments.

Kendall-Tau Correlation Experiment Result Tables. Figures 34a–37b provide Kendall-Tau experiment results for naturalness and token ratio correlations with linking performance (F1, Recall, and Precision).

- Figure 34a: Token-Character ratio–Recall
- Figures 35a–37b: Combined naturalness–[Recall, F1, Precision]
- Figures 38a–40b: [Regular, Low, Least] naturalness–recall
- Figures 41a–43b: [Regular, Low, Least] naturalness–F1
- Figures 44a–46b: [Regular, Low, Least] naturalness–precision
- Figures 47a–37b: [Combined, Regular, Low, Least] naturalness–execution accuracy

Database-level Schema Linking Box and Whisker Plots. Figures 51–54 show additional database-level box and whisker plots depicting schema linking performance over individual database schemas and their naturalness levels for each LLM.

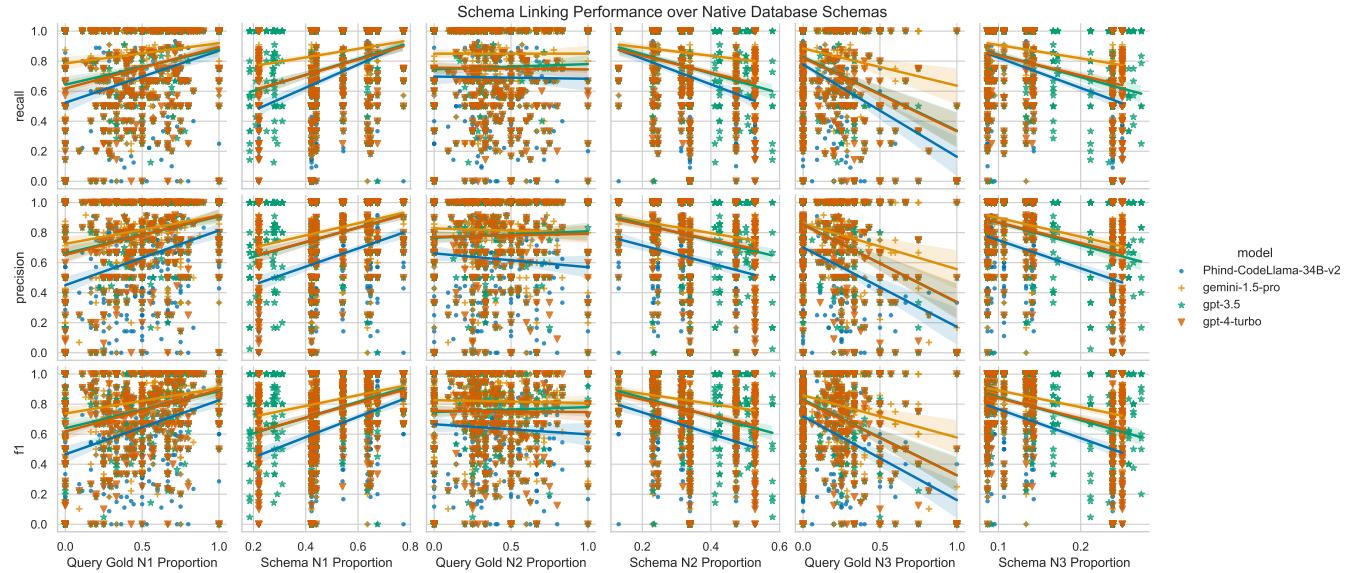


Figure 29: Scatterplot visualization of generated query schema linking performance in terms of F1, precision, and recall for identifier naturalness schema and gold query proportions in native schemas (modified schemas are excluded). The visualization indicates a correlation between naturalness levels and schema linking performance, where the presence of more natural identifiers corresponds to improved schema linking.

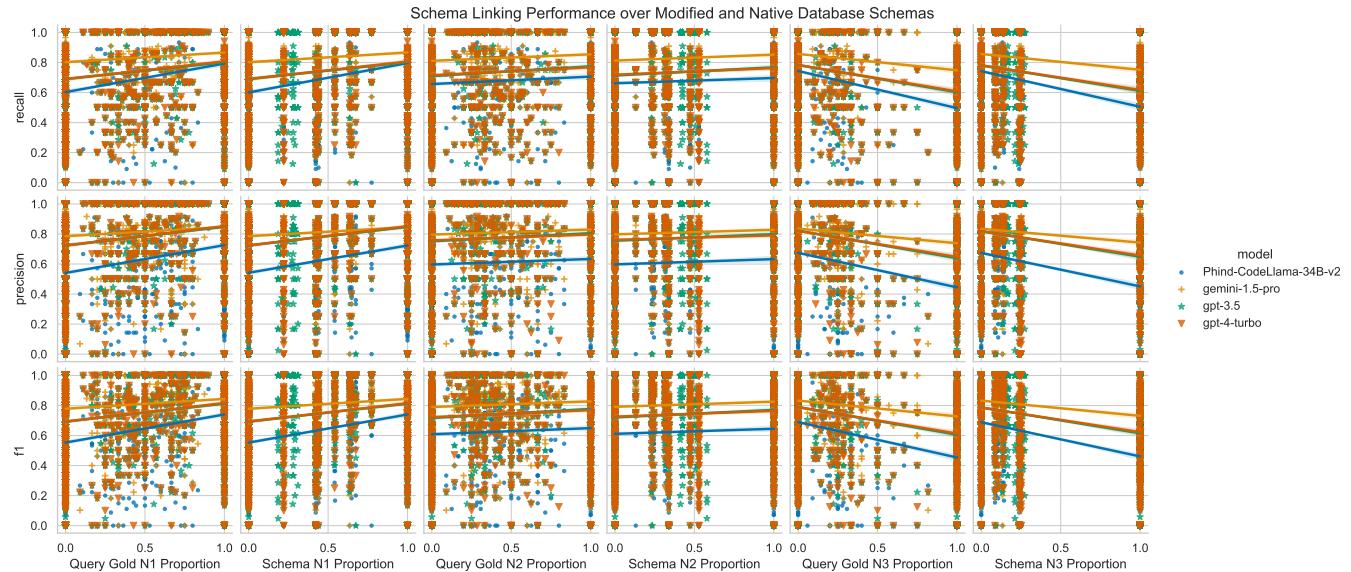


Figure 30: Scatterplot visualization of generated query schema linking performance in terms of F1, precision, and recall for identifier naturalness schema and gold query proportions in both modified and native schemas. The visualization indicates a correlation between naturalness levels and schema linking performance, where the presence of more natural identifiers corresponds to improved schema linking.

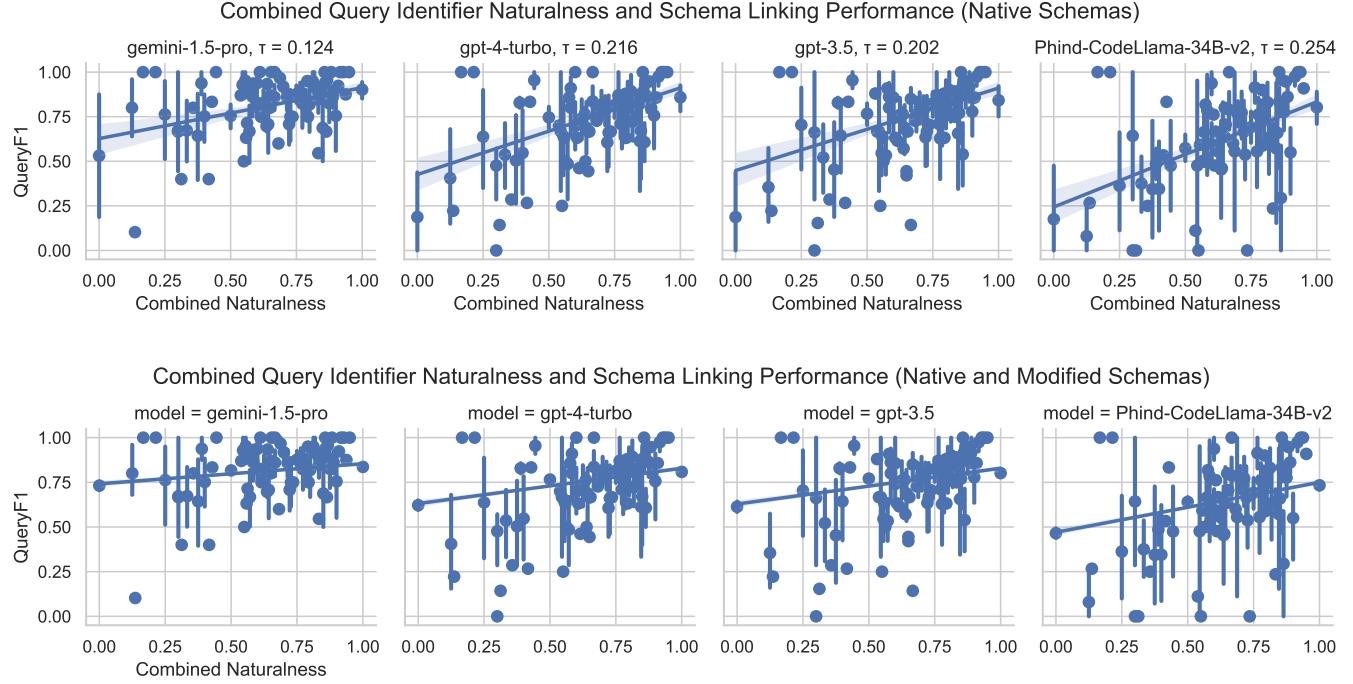


Figure 31: Schema linking performance (F1) of NL-SQL pairs relating to original schema naturalness levels. Regression line set with 0.95 confidence interval. Mean X value estimations for identical X values.

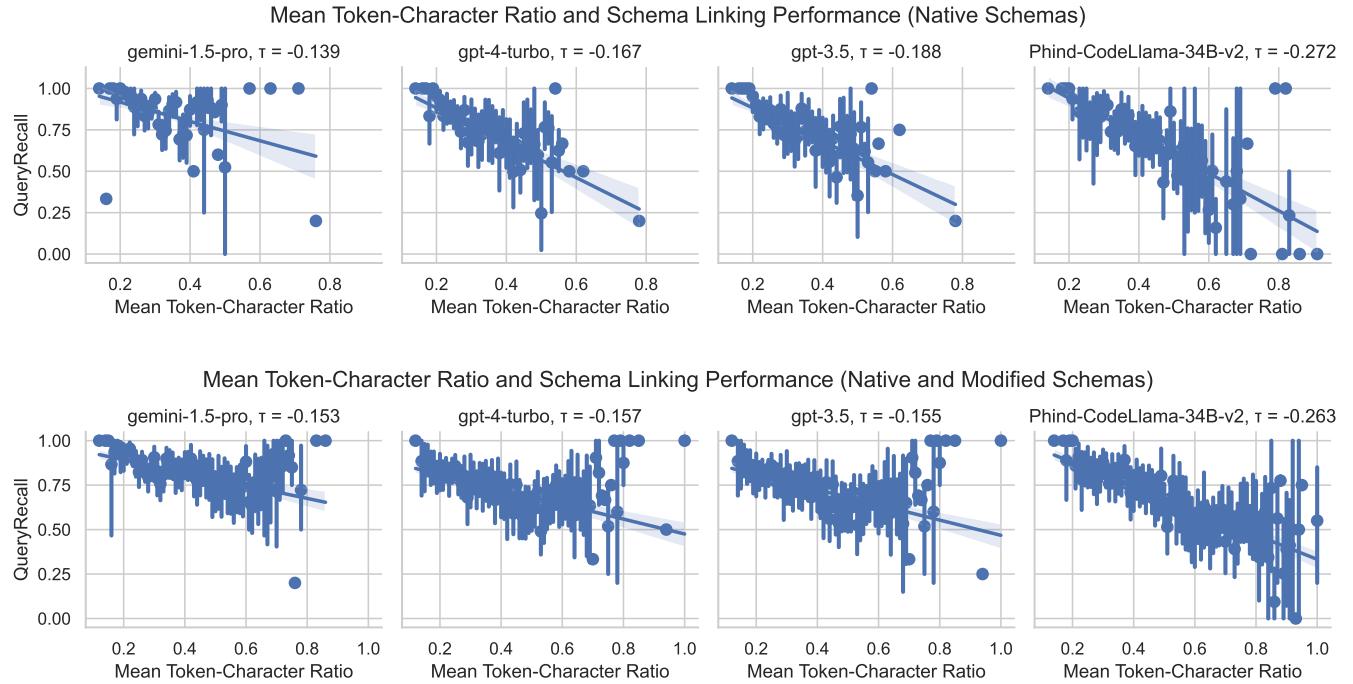


Figure 32: Schema linking performance (Recall) of NL-SQL pairs relating to query mean token-character ratios. Regression line set with 0.95 confidence interval. Mean X value estimations for identical X values.

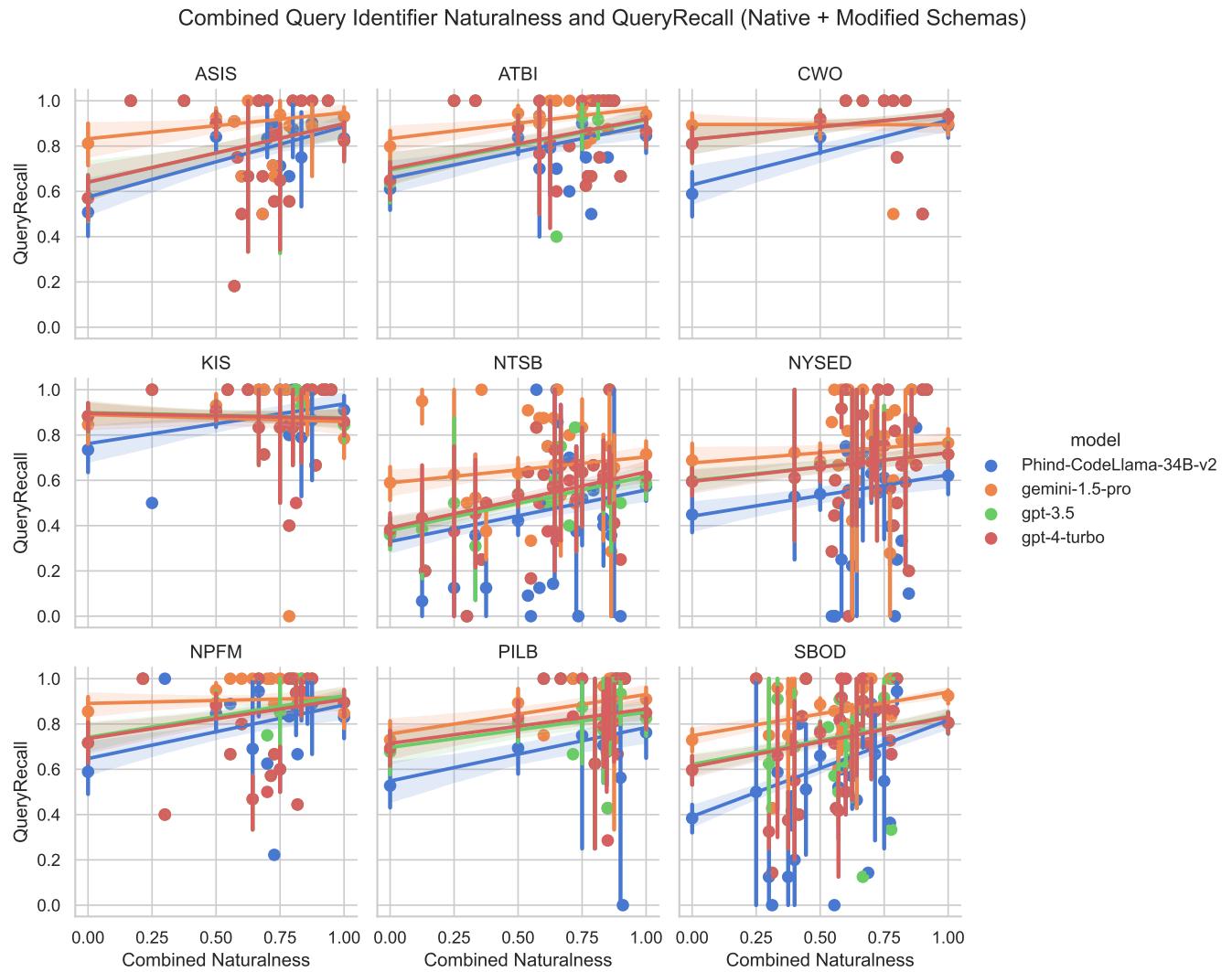


Figure 33: Database-level recall performance for modified and Native schema identifiers.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.159432	0.000002	492
gpt-4-turbo	-0.242731	0.000000	500
gpt-3.5	-0.234207	0.000000	500
Phind-CodeLlama-34B-v2	-0.272058	0.000000	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.143728	0.000000	1978
gpt-4-turbo	-0.166000	0.000000	1999
gpt-3.5	-0.166450	0.000000	1998
Phind-CodeLlama-34B-v2	-0.263041	0.000000	1936

(b) All schemas (native + modified)

Figure 34: Kendall-Tau Correlations between the mean query identifier token to character ratio and schema linking performance (Recall Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.137770	0.000066	492
gpt-4-turbo	0.223839	0.000000	500
gpt-3.5	0.209031	0.000000	500
Phind-CodeLlama-34B-v2	0.254438	0.000000	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.113625	0.000000	1978
gpt-4-turbo	0.176244	0.000000	1999
gpt-3.5	0.171700	0.000000	1998
Phind-CodeLlama-34B-v2	0.250113	0.000000	1936

(b) All schemas (native + modified)

Figure 35: Kendall-Tau Correlations between combined query identifier naturalness and schema linking performance (Recall Score) on native database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.151056	0.000005	492
gpt-4-turbo	0.221147	0.000000	500
gpt-3.5	0.214360	0.000000	500
Phind-CodeLlama-34B-v2	0.253227	0.000000	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.110138	0.000000	1978
gpt-4-turbo	0.185794	0.000000	1999
gpt-3.5	0.185058	0.000000	1998
Phind-CodeLlama-34B-v2	0.249083	0.000000	1936

(b) All schemas (native + modified)

Figure 36: Kendall-Tau Correlations between combined query identifier naturalness and schema linking performance (F1 Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.165627	0.000001	492
gpt-4-turbo	0.210652	0.000000	500
gpt-3.5	0.213515	0.000000	500
Phind-CodeLlama-34B-v2	0.240645	0.000000	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.099675	0.000000	1978
gpt-4-turbo	0.194037	0.000000	1999
gpt-3.5	0.193870	0.000000	1998
Phind-CodeLlama-34B-v2	0.240081	0.000000	1936

(b) All schemas (native + modified)

Figure 37: Kendall-Tau Correlations between combined query identifier naturalness and schema linking performance (Precision Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.117594	0.000720	492
gpt-4-turbo	0.203790	0.000000	500
gpt-3.5	0.179904	0.000000	500
Phind-CodeLlama-34B-v2	0.214458	0.000000	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.069195	0.000265	1978
gpt-4-turbo	0.129467	0.000000	1999
gpt-3.5	0.122796	0.000000	1998
Phind-CodeLlama-34B-v2	0.198193	0.000000	1936

(b) All schemas (native + modified)

Figure 38: Kendall-Tau Correlations between the proportion of query identifiers in the Regular naturalness category and schema linking performance (Recall Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.020559	0.555850	492
gpt-4-turbo	-0.046589	0.167214	500
gpt-3.5	-0.012805	0.704874	500
Phind-CodeLlama-34B-v2	-0.031016	0.359014	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.073173	0.000115	1978
gpt-4-turbo	0.062230	0.000708	1999
gpt-3.5	0.073867	0.000058	1998
Phind-CodeLlama-34B-v2	0.055309	0.002677	1936

(b) All schemas (native + modified)

Figure 39: Kendall-Tau Correlations between the proportion of query identifiers in the Low naturalness category and schema linking performance (Recall Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.165916	0.000010	492
gpt-4-turbo	-0.239739	0.000000	500
gpt-3.5	-0.239186	0.000000	500
Phind-CodeLlama-34B-v2	-0.285738	0.000000	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.158525	0.000000	1978
gpt-4-turbo	-0.213377	0.000000	1999
gpt-3.5	-0.212258	0.000000	1998
Phind-CodeLlama-34B-v2	-0.279290	0.000000	1936

(b) All schemas (native + modified)

Figure 40: Kendall-Tau Correlations between the proportion of query identifiers in the Least naturalness category and schema linking performance (Recall Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.131560	0.000074	492
gpt-4-turbo	0.198421	0.000000	500
gpt-3.5	0.182707	0.000000	500
Phind-CodeLlama-34B-v2	0.226842	0.000000	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.070395	0.000111	1978
gpt-4-turbo	0.141257	0.000000	1999
gpt-3.5	0.137642	0.000000	1998
Phind-CodeLlama-34B-v2	0.201929	0.000000	1936

(b) All schemas (native + modified)

Figure 41: Kendall-Tau Correlations between the proportion of query identifiers in the Regular naturalness category and schema linking performance (F1 Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.027036	0.417304	492
gpt-4-turbo	-0.042958	0.189129	500
gpt-3.5	-0.016076	0.623934	500
Phind-CodeLlama-34B-v2	-0.061160	0.059740	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.062756	0.000573	1978
gpt-4-turbo	0.053168	0.003026	1999
gpt-3.5	0.066247	0.000221	1998
Phind-CodeLlama-34B-v2	0.044565	0.012124	1936

(b) All schemas (native + modified)

Figure 42: Kendall-Tau Correlations between the proportion of query identifiers in the Low naturalness category and schema linking performance (F1 Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.158945	0.000009	492
gpt-4-turbo	-0.236171	0.000000	500
gpt-3.5	-0.243548	0.000000	500
Phind-CodeLlama-34B-v2	-0.240116	0.000000	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.147469	0.000000	1978
gpt-4-turbo	-0.217229	0.000000	1999
gpt-3.5	-0.220154	0.000000	1998
Phind-CodeLlama-34B-v2	-0.263193	0.000000	1936

(b) All schemas (native + modified)

Figure 43: Kendall-Tau Correlations between the proportion of query identifiers in the Least naturalness category and schema linking performance (F1 Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.140165	0.000039	492
gpt-4-turbo	0.178184	0.000000	500
gpt-3.5	0.177250	0.000000	500
Phind-CodeLlama-34B-v2	0.220023	0.000000	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.063466	0.000685	1978
gpt-4-turbo	0.153237	0.000000	1999
gpt-3.5	0.148501	0.000000	1998
Phind-CodeLlama-34B-v2	0.196175	0.000000	1936

(b) All schemas (native + modified)

Figure 44: Kendall-Tau Correlations between the proportion of query identifiers in the Regular naturalness category and schema linking performance (Precision Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.025641	0.453269	492
gpt-4-turbo	-0.005640	0.867278	500
gpt-3.5	0.003413	0.919671	500
Phind-CodeLlama-34B-v2	-0.060758	0.065561	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.054467	0.003577	1978
gpt-4-turbo	0.041575	0.024761	1999
gpt-3.5	0.056122	0.002463	1998
Phind-CodeLlama-34B-v2	0.041782	0.020555	1936

(b) All schemas (native + modified)

Figure 45: Kendall-Tau Correlations between the proportion of query identifiers in the Low naturalness category and schema linking performance (Precision Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.170160	0.000004	492
gpt-4-turbo	-0.237864	0.000000	500
gpt-3.5	-0.248686	0.000000	500
Phind-CodeLlama-34B-v2	-0.215408	0.000000	484

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.135749	0.000000	1978
gpt-4-turbo	-0.215227	0.000000	1999
gpt-3.5	-0.222732	0.000000	1998
Phind-CodeLlama-34B-v2	-0.247999	0.000000	1936

(b) All schemas (native + modified)

Figure 46: Kendall-Tau Correlations between the proportion of query identifiers in the Least naturalness category and schema linking performance (Precision Score) on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.031553	0.402105	503
gpt-4-turbo	0.127587	0.000704	503
gpt-3.5	0.130571	0.000526	503
Phind-CodeLlama-34B-v2	0.133928	0.000376	503

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.051044	0.013769	2012
gpt-4-turbo	0.105441	0.000000	2012
gpt-3.5	0.113802	0.000000	2012
Phind-CodeLlama-34B-v2	0.156595	0.000000	2012

(b) All schemas (native + modified)

Figure 47: Kendall-Tau Correlations between the proportion of query identifiers in the Regular naturalness category and Execution Accuracy on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.067135	0.075780	503
gpt-4-turbo	-0.028380	0.452870	503
gpt-3.5	-0.013996	0.711227	503
Phind-CodeLlama-34B-v2	-0.051436	0.173713	503

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.065896	0.001482	2012
gpt-4-turbo	0.035483	0.087019	2012
gpt-3.5	0.059034	0.004410	2012
Phind-CodeLlama-34B-v2	0.025475	0.219203	2012

(b) All schemas (native + modified)

Figure 48: Kendall-Tau Correlations between the proportion of query identifiers in the Low naturalness category and Execution Accuracy on native and modified database schemas.

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.170360	0.000027	503
gpt-4-turbo	-0.198597	0.000001	503
gpt-3.5	-0.209644	0.000000	503
Phind-CodeLlama-34B-v2	-0.167881	0.000035	503

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	-0.153130	0.000000	2012
gpt-4-turbo	-0.178487	0.000000	2012
gpt-3.5	-0.193995	0.000000	2012
Phind-CodeLlama-34B-v2	-0.210597	0.000000	2012

(b) All schemas (native + modified)

Figure 49: Kendall-Tau Correlations between the proportion of query identifiers in the Least naturalness category and Execution Accuracy on native and modified database schemas.

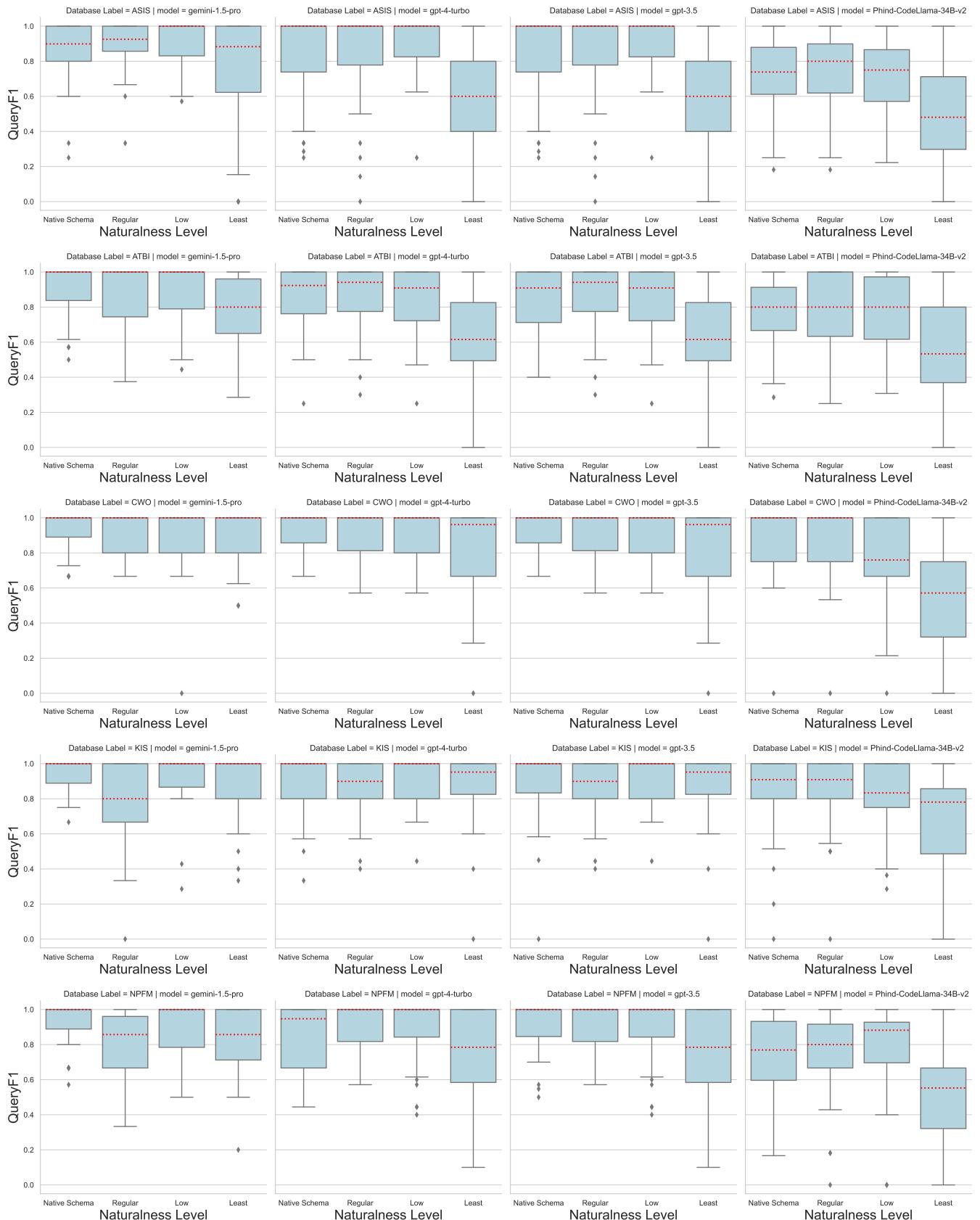
Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.077353	0.038579	503
gpt-4-turbo	0.153124	0.000042	503
gpt-3.5	0.161218	0.000016	503
Phind-CodeLlama-34B-v2	0.151421	0.000051	503

(a) Native schemas

Model	Kendall-Tau	P Value	n
gemini-1.5-pro	0.095640	0.000002	2012
gpt-4-turbo	0.143250	0.000000	2012
gpt-3.5	0.156622	0.000000	2012
Phind-CodeLlama-34B-v2	0.195768	0.000000	2012

(b) All schemas (native + modified)

Figure 50: Kendall-Tau Correlations between combined query identifier naturalness and execution accuracy on native and modified database schemas.

**Figure 51: Schema linking performance (F1 score) changes across database naturalness levels.**

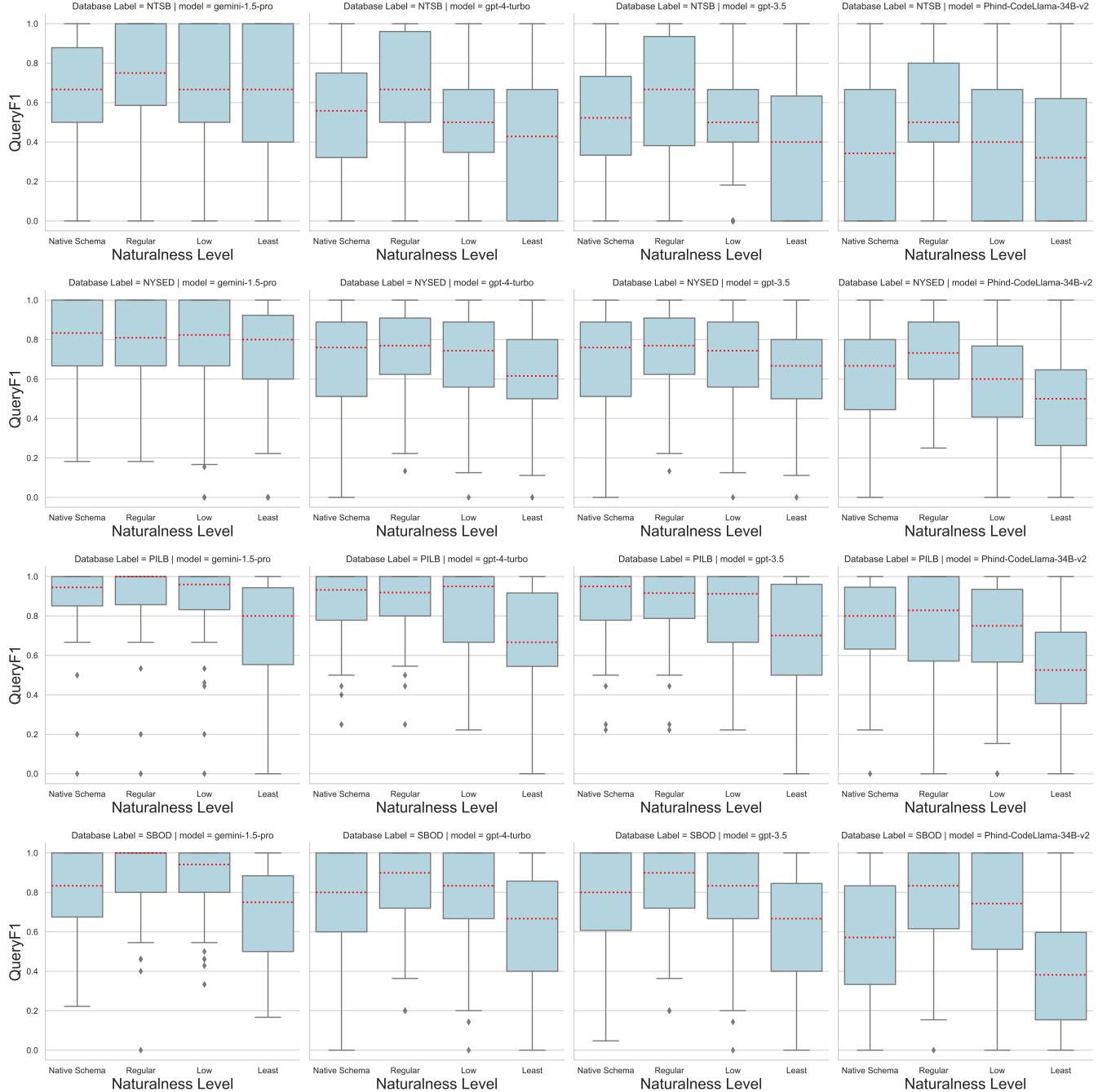


Figure 52: Schema linking performance (F1 score) changes across database naturalness levels.

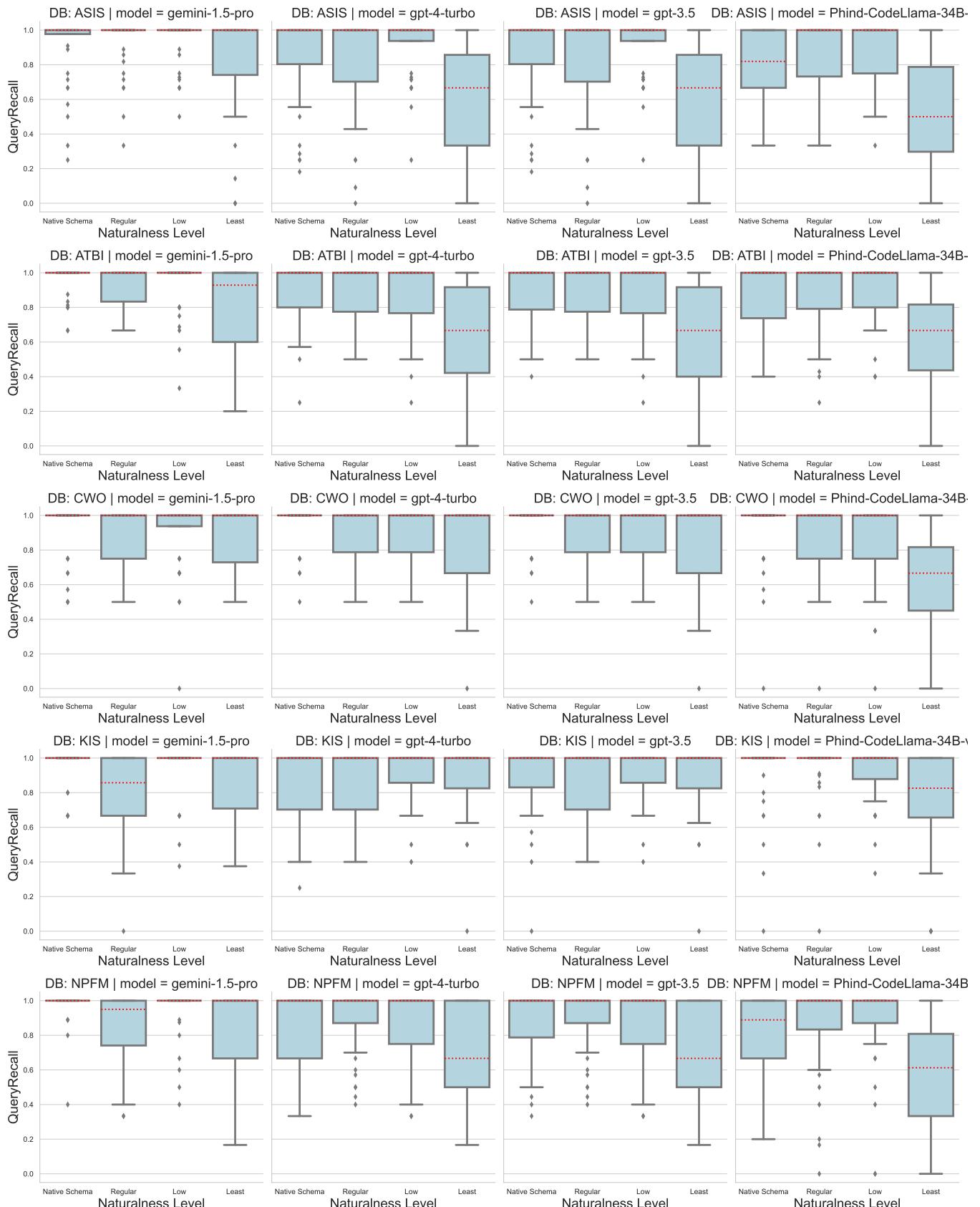


Figure 53: Schema linking performance (Recall score) changes across database naturalness levels.

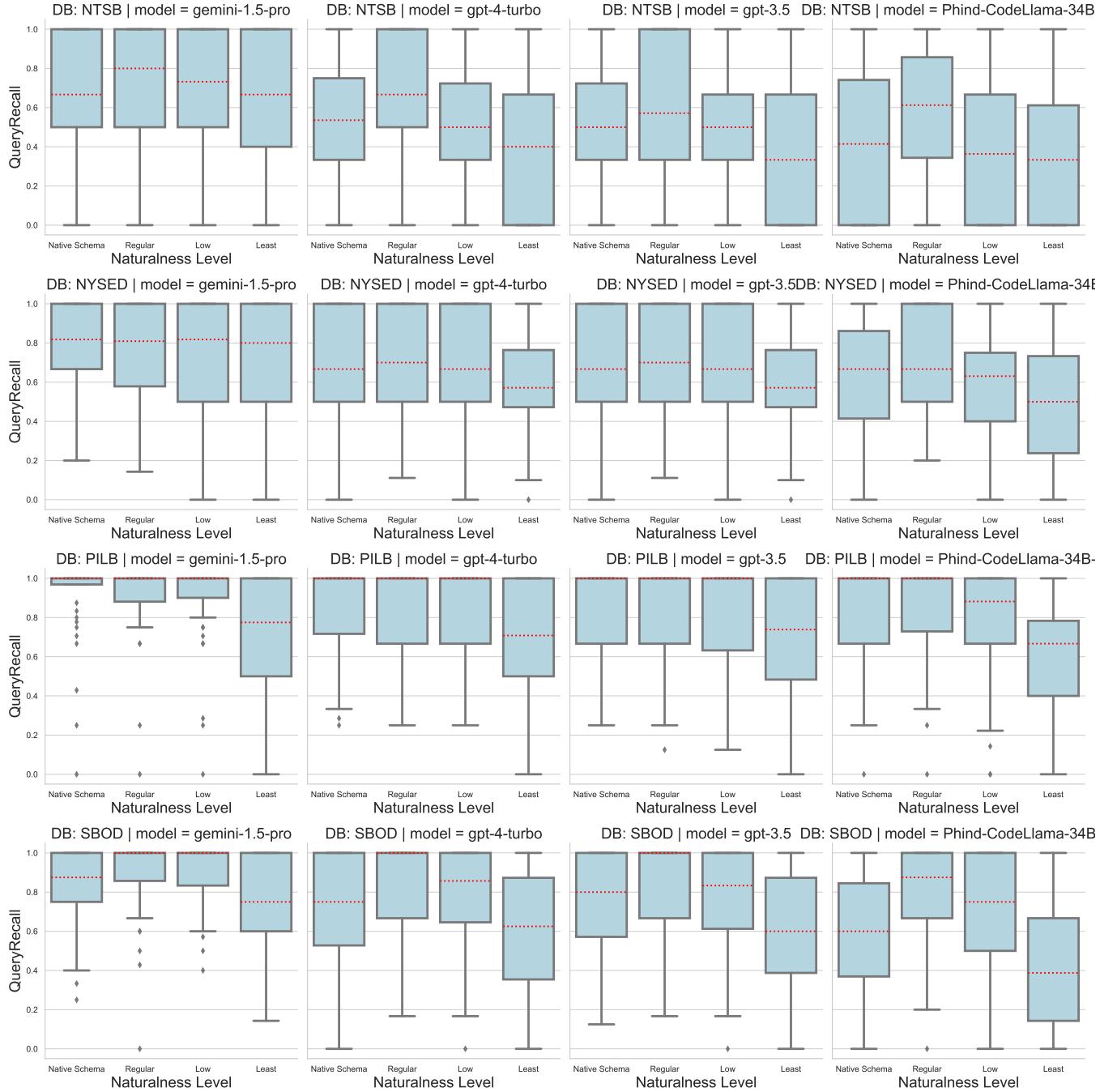


Figure 54: Schema linking performance (Recall score) changes across database naturalness levels.