

ArtifactLinker: Linking Scientific Artifacts for Automatic State-of-the-art Discovery

Anonymous ACL submission

Abstract

Scientific artifacts such as models and benchmarks underpin machine learning research. With the rapid growth of repositories like HuggingFace, researchers now have access to millions of artifacts, yet a key challenge remains: how can we automatically discover the state-of-the-art (SOTA) model for a given benchmark by fully leveraging existing artifacts? We formalize this as **automatic SOTA discovery** by modeling HuggingFace as an artifact graph, where nodes are models/benchmarks and edges represent evaluations. We propose ARTIFACTLINKER, a two-stage framework: (1) prediction of promising unobserved model–benchmark links using Graph Neural Networks (GNNs) or graph-augmented Large Language Models (LLMs), and (2) verification via fully automatic, reproducible coding experiments with agents. We further introduce ARTIFACTBENCH with 2,977 models and 559 benchmarks to evaluate for both stages. Results show effective graph-based prediction and reliable end-to-end automatic verification of high-performing candidates.

1 Introduction

Scientific artifacts are the fundamental building blocks of research (Heumüller et al., 2020; Cooper et al., 2022; Johnson et al., 2019). Codebases on the GitHub platform, papers available on arXiv, models and benchmarks on the HuggingFace are all examples of such artifacts. Researchers engaged in doing reproducible and high-quality research share, interact with, and build upon these artifacts, releasing new versions to demonstrate progress (Marić et al., 2023; Lissa et al., 2020). In the machine learning community, a vast number of artifacts (>1M on HuggingFace) are produced by researchers working in different domains (Castaño et al., 2024; Ait et al., 2023; Laufer et al., 2025). This naturally raises an important question: *How can we lever-*

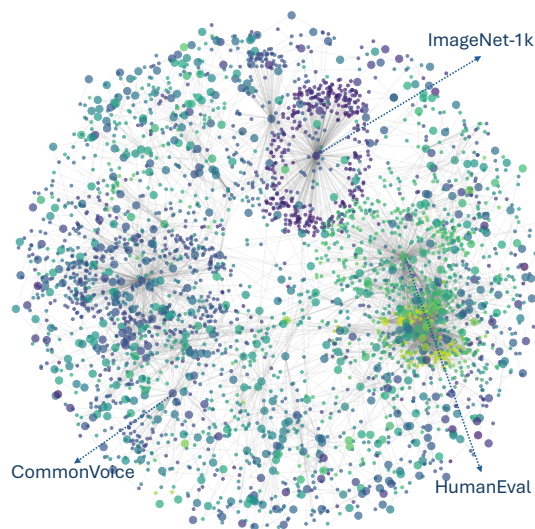


Figure 1: This shows the underlying graph of research artifacts from Huggingface Hub, used to build ARTIFACTBENCH, a new suite of tasks for developing automatic discovery agents. Large/small nodes represent models/datasets linked through evaluation, which we leverage for a new task called **SOTA discovery** (i.e., discovering novel links between existing artifacts).

age existing artifacts to enable automatic discovery? Addressing this question would (1) allow us to utilize diverse types of artifacts better, and (2) promote scalable and automated scientific discovery based on existing resources. We focus on the HuggingFace community as a case study, since it is one of the largest and most active hubs of open-source machine learning artifacts and aims to make experiments more accessible and easy to run. With countless models, benchmarks, and libraries hosted on the platform, it provides an invaluable foundation for exploring automated discovery.

Building an automatic discovery system on HuggingFace presents several challenges. First, the concept of “automatic discovery” itself is ill-defined (Beel et al., 2025; Kitano, 2021; Kramer et al., 2023)—what does it really mean for a system to conduct research autonomously and con-

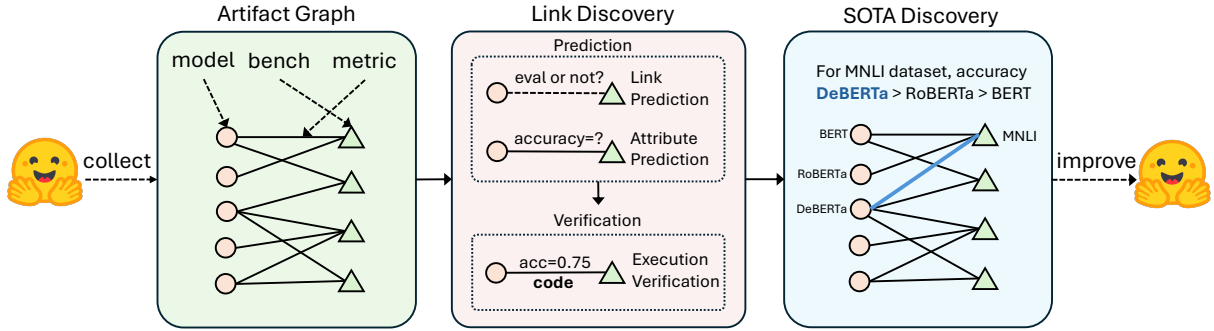


Figure 2: **Overview of ARTIFACTLINKER.** The overall automatic SOTA discovery process in ARTIFACTLINKER includes three main components: (1) artifact graph construction from HuggingFace (models, benchmarks, evaluation relationships); (2) link prediction and verification as discovery tasks; (3) state-of-the-art result filtering.

060 tribute to future applications? Second, although
 061 HuggingFace provides convenient access to mod-
 062 els and datasets, building a fully automated and
 063 reproducible pipeline to verify predicted model per-
 064 formance is still difficult (Urbanowicz et al., 2022).
 065 The usability of models and benchmarks remains
 066 inconsistent, and many of them require special-
 067 ized configurations to work properly and reproduce
 068 research works built on them, even a frustrating
 069 process for human researchers (Mu et al., 2025).

070 **Thinking of HuggingFace as an artifact graph**

071 Our key observation is that the HuggingFace
 072 community can be naturally represented as a
 073 graph (Chen et al., 2025; Laufer et al., 2025) as
 074 shown in Figure 1 and Figure 2 (left), where mod-
 075 els and benchmarks serve as nodes, and their rela-
 076 tionships form the edges. This perspective is moti-
 077 vated by three characteristics of the platform: (1) it
 078 hosts a large collection of artifacts, including mod-
 079 els and benchmarks; (2) many artifacts are high-
 080 quality and widely adopted; and (3) it encodes rich
 081 relational information—for example, a model trained
 082 on one dataset and evaluated on another with a re-
 083 ported F1 score. Such structured relationships are
 084 often difficult to extract directly from research pa-
 085 pers. Prior work has largely treated HuggingFace
 086 as an information source for retrieval (Silva et al.,
 087 2025) or as an API hub (Shen et al., 2023). In
 088 contrast, we highlight its value for dynamic dis-
 089 covery. Rather than viewing it as a static data
 090 source, we aim to actively uncover new relation-
 091 ships—particularly model–benchmark interactions
 092 defined through evaluation metrics. These interac-
 093 tions naturally form the edges of an artifact graph.

094 **Linking artifacts as automatic SOTA discovery**

095 Within this artifact graph, we define the task of
 096 automatic discovery as finding links with state-of-
 097 the-art evaluation results between artifacts (models
 098 and benchmarks). Concretely, our goal is to build

an auto-discovery engine that leverages the existing
 artifact graph structures in the HuggingFace com-
 munity and their existing relationships to propose
 new promising links.

Novel framework for automatic discovery To
 operationalize the linking problem as automatic
 discovery, we propose a novel framework with two
 stages (see Figure 2 (center)): (1) prediction and
 (2) verification. Given the vast number of artifacts
 in the graph and the high cost of verifying evalu-
 ation results, directly identifying the small subset
 of valuable links through execution alone is pro-
 hibitively costly. Our framework addresses this
 challenge by first performing prediction, where
 we use strong priors to filter out the majority of
 unlikely links—analogueous to how experienced re-
 searchers prioritize promising directions. In the
 verification stage, we employ a coding agent to test
 and validate the predicted links, grounding hypoth-
 eses into real, reproducible results. This division of
 labor enables scalable automatic discovery while
 ensuring that results remain empirically verifiable.

Contributions Our work makes three key con-
 tributions: (1) we construct ARTIFACTBENCH, a
 challenging new benchmark for prediction, veri-
 fication, and automatic discovery, (2) we design
 a two-stage framework named ARTIFACTLINKER
 for scalable link discovery, and (3) we demonstrate
 the practical value of this system through extensive
 evaluation and highlight current limitations that
 motivate further research on ARTIFACTBENCH.

099 **2 Related Work**

HuggingFace platform utilization Hugging-
 Face has increasingly become a natural platform
 for studying automatic discovery. Prior work has
 largely relied on static analyses of its artifacts
 and relationships to characterize trends in machine
 learning development (Chen et al., 2025; Laufer

et al., 2025). Beyond serving as a repository, HuggingFace has been conceptualized in multiple ways: as a knowledge graph (Silva et al., 2025), an API hub (Shen et al., 2023), a model card aggregator (Yang et al., 2024), and even an evolutionary tree (Gao and Gao, 2023). Other studies have examined its community dynamics (Rahman et al., 2025; Castaño et al., 2023). In contrast, our work moves beyond static description and trend analysis and focuses on execution-based prediction/verification.

Large-scale prediction for accelerating discoveries Accelerating scientific discovery has been a major focus in domains such as drug discovery (Stokes et al., 2020; Serrano et al., 2024; Vişan and Neguţ, 2024; You et al., 2022), materials science (Xie and Grossman, 2018; Butler et al., 2018), and molecular design (Segler et al., 2018), among others (Cheng et al., 2025). In these settings, experimental verification is prohibitively costly and time-consuming. In contrast, our work focuses on a more tractable class of automatic discovery tasks by leveraging the intrinsic linking structure of HuggingFace artifacts.

LLM-based coding agents for reproducible experimentation Prior work has explored free-form discovery with generating executable code from research ideas (Lu et al., 2024; Jansen et al., 2024, 2025), though evaluation remains challenging given the open-ended nature of such tasks. Other efforts have focused on reproducing experiments within specific codebases (Bogin et al., 2024; Starace et al., 2025; Kim et al., 2025; Seo et al., 2025; Siegel et al., 2024; Xiang et al., 2025), which is challenging due to the complexity of such codebases. In contrast, our tasks rely on reproducing a more concrete/grounded set of research artifacts.

3 Constructing an Artifact Graph from HuggingFace Hub

As the first stage of ARTIFACTLINKER, we first formally provide the definition of artifact graphs that we conduct link discovery on. Furthermore, we provide details about how we extract the artifact graph from the HuggingFace platform.

Definition of artifact graphs We define the artifact graph as an undirected heterogeneous bipartite graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where the node set $\mathcal{V} = \mathcal{V}_m \cup \mathcal{V}_d$ consists of two disjoint types of nodes: **model nodes** \mathcal{V}_m and **benchmark nodes** \mathcal{V}_d . The edge set $\mathcal{E} \subseteq \mathcal{V}_m \times \mathcal{V}_d \times \mathcal{K}$ encodes **evaluation relationships**, where each edge (u, v, k) indicates that a

model $u \in \mathcal{V}_m$ has an evaluation result on a benchmark $v \in \mathcal{V}_d$ with score $k = \phi(u, v)$ defined by a metric function $\phi : \mathcal{V}_m \times \mathcal{V}_d \rightarrow \mathcal{K}$ (e.g., accuracy or F-score). Besides edges, each benchmark node is associated with attributes such as metadata (e.g., download counts) and task descriptions, while each model node is associated with metadata and specifications, including architecture, number of parameters, and configuration. Such rich information on both edges and nodes allows the graph to capture fine-grained performance relationships between models and benchmarks.

Data collection We construct edges in the artifact graph by extracting ground-truth evaluation metrics from HuggingFace model cards. Specifically, we parse the README files to identify model–dataset pairs along with their reported performance scores. Model and dataset names are then matched to their canonical entries on the HuggingFace platform to ensure consistency, yielding a clean set of evaluation edges. In total, this process produces $|\mathcal{E}| = 5,476$ perfectly matched evaluation records, which serve as edge attributes in the graph. For nodes, we collect both models and benchmarks. Each model is described by its model card, dataset card, and metadata fields (e.g., architecture, parameters, tags), while each dataset includes metadata such as domain, size, and license. To ensure relevance, we retain only nodes that participate in at least one evaluation edge (isolated artifacts).

Graph statistics The final artifact graph contains $|\mathcal{V}| = 3,536$ nodes, consisting of $|\mathcal{V}_m| = 2,977$ models and $|\mathcal{V}_d| = 559$ benchmarks, with $|\mathcal{E}| = 5,476$ edges in total. To quantify the community structure, we apply Louvain community detection (Blondel et al., 2008) and compute the modularity, obtaining $Q = 0.82$, which highlights the presence of clear sub-community structures under sparse connectivity. Figure 1 shows the visualization of the artifact graph structure. While benchmarks like ImageNet (Deng et al., 2009) and HumanEval (Chen, 2021) form expected hubs, such analysis indicates that edge concentration is limited and not solely focused on high-degree hubs.

4 Linking Scientific Artifacts for Automatic SOTA Discovery

To describe the overall pipeline of ARTIFACTLINKER, we first formalize the problem definition of automatic discovery using the artifact-graph formulation. We then introduce our scalable

solution, which addresses this problem via a two-stage *prediction–verification* framework.

4.1 Definition of Automatic SOTA Discovery

Building on the artifact graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ defined in the previous section, each observed edge $(m, d) \in \mathcal{E}$ is annotated with a performance score given by the metric function $\phi : \mathcal{V}_m \times \mathcal{V}_d \rightarrow \mathcal{K}$, with a realized score denoted as $\phi(m, d)$. The ultimate goal of automatic discovery is to identify a missing link $(m, d, k) \notin \mathcal{E}, m \in \mathcal{V}_m, d \in \mathcal{V}_d$ such that the predicted score under the same evaluation function $\phi(m, d)$ exceeds the best known performance on dataset d :

$$\phi(m, d) > \max_{(m', d) \in \mathcal{E}} \phi(m', d).$$

In other words, automatic SOTA discovery seeks model–benchmark pairs not yet connected in \mathcal{G} but expected to advance the state of the art on d . The input of our proposed ARTIFACTLINKER framework is the artifact graph \mathcal{G} while the target output is these state-of-the-art (m, d) pairs.

4.2 Scalable Link Discovery Framework

Because explicit model evaluation is computationally expensive, scalability requires pruning candidates before verification. We therefore adopt a two-stage framework with prediction and verification. In the prediction stage, we rank model–dataset pairs (m, d) using a score function $s(m, d)$ and select a candidate set \mathcal{C} whose predicted scores approach or exceed the current SOTA. In the verification stage, candidates in \mathcal{C} are executed and evaluated, producing the validated SOTA set \mathcal{C}^* .

4.2.1 Predicting via Score Function Modeling

The goal of prediction is not only to estimate $\hat{\phi}(m, d)$ for missing edges but also to rank all candidate links by their discovery potential. Concretely, each prediction model defines a scoring function $s : \mathcal{V}_m \times \mathcal{V}_d \rightarrow \mathbb{R}$, which assigns each unseen pair $(m, d) \notin \mathcal{E}$ a predicted performance. For each model the objective is to approximate the score function $\hat{\phi}(m, d)$ and produce rankings consistent with observed edges. Since all of our prediction models utilize graph information, we formulate them using a general GNN-style notation. Formally, each score function is defined as:

$$\mathbf{h}_v^{(k+1)} = \text{AGG}^{(k)}(\mathbf{h}_v^{(k)}, \{\mathbf{h}_u^{(k)}, e_{uv} : u \in \mathcal{N}(v)\})$$

$$s(m, d) = \psi(\mathbf{h}_m^{(L)}, \mathbf{h}_d^{(L)})$$

where $\mathbf{h}_v^{(k)}$ is a model’s internal representation of node v after k iterations aggregated over all representations $\mathbf{h}_u^{(k)}$ constructed from neighboring nodes $u \in \mathcal{N}(v)$. The final score is $s(m, d)$ is then defined over the final node representations $\mathbf{h}_m^{(L)}$ and $\mathbf{h}_d^{(L)}$ transformed via some function ψ . As we detail below, different instantiations of AGG and ψ lead to different modeling strategies.

LLM modeling These models are based on ordinary prompt-based LLMs. Following the TextGNN formulation presented in Yu et al. (2024), we represent each $\mathbf{h}_v^{(k)}$ directly with text instead of embeddings. The aggregation function AGG is instantiated as a one-layer *concat* operator that serializes neighborhood information into a textual prompt:

$$\mathbf{h}_v^{(k+1)} = \text{CONCAT}(\mathbf{h}_v^{(k)}, \{\mathbf{h}_u^{(k)}, e_{uv} : u \in \mathcal{N}(v)\})$$

$$s(m, d) = \text{LLM}(\mathbf{h}_m^{(L)}, \mathbf{h}_d^{(L)})$$

Here, $\mathbf{h}_v^{(0)} = \mathbf{x}_v$ corresponds to the textual description of node v . The CONCAT function converts all neighbor attributes into natural language and concatenates them into a serialized context. The LLM then consumes these prompts for m and d , jointly reasoning over their neighborhoods to predict $s(m, d)$. We typically adopt a single TextGNN layer due to the limited context window size.

GNN-based modeling Alternatively, we instantiate AGG using standard message passing in Graph Neural Networks. In this setting, node features are initialized as $\mathbf{h}_v^{(0)} = \text{EMB}(\mathbf{x}_v)$ from the textual description \mathbf{x}_v , and representations are updated by:

$$\mathbf{h}_v^{(k+1)} = \sigma(W^{(k)}\mathbf{h}_v^{(k)} + \text{AGG}^{(k)}(\mathcal{N}(v)))$$

$$s(m, d) = \text{MLP}([\mathbf{h}_m^{(L)} \parallel \mathbf{h}_d^{(L)}])$$

Here, $\text{AGG}^{(k)}(\cdot)$ can be implemented with classical GNN operators such as GATv2Conv (Brody et al., 2021), which applies dynamic attention weights to neighbor embeddings. The edge scoring function $s(m, d)$ is then computed by an MLP over the concatenation of the final node embeddings $\mathbf{h}_m^{(L)}$ and $\mathbf{h}_d^{(L)}$. In contrast to the LLM modeling above, this model is tuned using a training set of positive and negative links mined from the original artifact graph (see full details in Appendix C.1)

Ranking with score function Given different modeling choices of $s(m, d)$, the prediction stage first produces a ranked list of all unseen pairs $(m, d) \notin \mathcal{E}$, ordered by their predicted score

$s(m, d)$. This ranking reflects the discovery potential of each candidate, with higher scores indicating stronger likelihood of advancing the state of the art. From this ranked list, we further define the set of promising candidates as

$$\mathcal{C} = \left\{ (m, d) \left| \begin{array}{l} (m, d) \notin \mathcal{E}, m \in \mathcal{V}_m, d \in \mathcal{V}_d, \\ s(m, d) \geq \max_{(m', d) \in \mathcal{E}} (1 - \delta) \phi(m', d) \end{array} \right. \right\}$$

where $\delta(d) \in [0, 1]$ is a dataset-specific tolerance parameter and $\phi(m', d)$ is the existing evaluation score of model m' on dataset d under metric $\mu(d)$. Larger $\delta(d)$ indicates higher verification cost, but the greater possibility for finding the state-of-the-art performance. In other words, the score function $s(m, d)$ serves a dual purpose: it provides a global ranking over all candidates, and it defines a thresholded subset \mathcal{C} to guide the verification stage by focusing on the most promising discoveries.

4.2.2 Verification via Multi-turn Tool Use

Given the candidate set \mathcal{C} of manageable size, the goal of the verification stage is to execute each model–dataset pair and obtain concrete evaluation results. For each candidate (m, d) , we perform verification using the CodeAct agentic workflow (Wang et al., 2024), which plans, generates, executes, and iteratively refines code in a multi-turn manner for scalable automatic benchmarking.

Specifically, the agent takes as input the model identifier m , the dataset identifier d , and a target evaluation metric $\mu(d)$ hooked with the dataset d (e.g., accuracy, F1, BLEU). We represent these inputs as descriptors m , d , and μ . To support reliable execution, the agent is equipped with HuggingFace-specific tools—such as `web_search`, `get_metadata`, and `get_readme`—which can be invoked at each turn to retrieve relevant model and dataset information. This tool-augmented agentic workflow enables reliable automatic evaluation:

$$s_{\mu}^*(m, d) = \text{CODEACT}(m, d, \mu(d)).$$

4.3 ARTIFACTLINKER Algorithm

Algorithm 1 presents the pipeline of ARTIFACTLINKER for automatic SOTA discovery. The algorithm proceeds in three stages: (i) compute the current SOTA (ii) perform attribute prediction over a fixed artifact graph (a transductive setting where nodes are observed and only missing links are predicted), with GNN-based methods enabling batch-parallel inference; and (iii) verify top candidates with a coding agent. Larger δ admits more

Algorithm 1: SOTA discovery workflow

Input: Artifact graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$; score model s ; tolerance δ

Output: Validated SOTA discoveries \mathcal{C}^*

Precompute current SOTA (per dataset);

foreach $d \in \mathcal{V}_d$ **do**

$b(d) \leftarrow \max_{(m', d) \in \mathcal{E}} \phi(m', d)$

Prediction: propose promising links;

$\hat{s} \leftarrow \text{ATTRPREDICT}(s, \mathcal{G})$;

$\mathcal{C} \leftarrow \{(m, d) \notin \mathcal{E} \mid \hat{s}(m, d) \geq (1 - \delta)b(d)\}$;

Verification: execute and validate;

$\mathcal{C}^* \leftarrow \emptyset$;

foreach $(m, d) \in \mathcal{C}$ **do**

$s_{\mu}^*(m, d) \leftarrow \text{CODEACT}(m, d, \mu(d))$

if $s_{\mu}^*(m, d) > b(d)$ **then**

$\mathcal{C}^* \leftarrow \mathcal{C}^* \cup \{(m, d)\}$;

return \mathcal{C}^* ;

candidates and thus increases computational overhead during verification. As a special case, ARTIFACTLINKER can also operate in a conditional mode, restricted to a specific dataset d or model m , reducing the need for full graph traversal. In an inductive setting, where a new model or benchmark is introduced, ARTIFACTLINKER can incorporate the new node via its attributes and apply the same predict-and-verify procedure to propose evaluations involving the unseen artifact.

5 Evaluating Automatic SOTA Discovery

To evaluate ARTIFACTLINKER, we design a mask-and-predict/verify protocol based on the artifact graph. Instead of running in the real world without ground-truths, we directly sample existing links $(m, d) \in \mathcal{E}$ and evaluate them on ARTIFACTBENCH by edge masking.

For evaluation, we consider two forms of missing links/masking. (i) *Edge masking*: a subset of observed edges $\mathcal{E}_{\text{mask}} \subset \mathcal{E}$ is hidden from the graph, requiring the model to recover their existence. (ii) *Attribute masking*: for another subset of edges, we keep the edge structure visible but hide their associated attributes (e.g., evaluation scores), requiring the model to predict/rank the missing values.

Prediction evaluation Given the masked graph $\mathcal{G} \setminus \mathcal{E}_{\text{mask}}$, we benchmark prediction performance with four tasks: (1) *Link prediction (edge masking)*: determine whether a masked edge $(m, d) \in \mathcal{E}_{\text{mask}}$

exists, i.e., $\hat{y}_{m,d} = \mathbf{1}[s(m,d) > \tau]$, where $s(m,d)$ is the predicted score function. (2) *Link ranking (edge masking)*: for each dataset d , rank all candidate models m by $s(m,d)$ and evaluate the rank of the true masked edges $(m,d) \in \mathcal{E}_{\text{mask}}$. (3) *Attribute prediction (attribute masking)*: for each edge (m,d) with hidden attributes, predict its score $\hat{\phi}(m,d)$. (4) *Attribute ranking (attribute masking)*: for each dataset d , rank all candidate models m by predicted attributes $\hat{\phi}(m,d)$ and compare this ordering against ground-truth scores.

Verification evaluation Finally, we define the *reproduction task* in the verification stage, which also corresponds to attribute masking. Here, the coding agent re-evaluates edges (m,d) with hidden attributes to measure consistency between predicted scores $s(m,d)$ and execution-based results $\phi^*(m,d)$. Formally, consistency is satisfied if $|s(m,d) - \phi^*(m,d)| \leq (1 - \delta)\phi^*(m,d)$. We choose $\delta=0.1$. This criterion checks whether prediction and execution agree up to an acceptable margin, ensuring faithfulness of the verification.

6 Experimental Settings

Prediction task settings We perform transductive edge-level splitting on the artifact graph, randomly partitioning edges into training, development, and test sets with a 70%–10%–20% ratio. The node set is fixed and shared across all splits, as our goal is to predict missing links in an existing graph. During training, message passing is restricted to the training edges, while development and test edges are held out for evaluation.

Prediction baseline settings Besides the LLM-based TextGNN settings and GNN-based settings, we include multiple baseline settings for comparison: (1) *metadata*, where we directly use metadata of the artifacts such as downloading times and metrics on edges as a feature for prediction; (2) *graph baselines*, where we utilize graph-related algorithm including Matrix Factorization (MF), Adamic-Adar (AA), Common-Neighborhood (CN) and more modern NeoGNN (Yun et al., 2021) baselines; (4) *LLM baselines*, where we only provide the model and benchmark information to an LLM (we use GPT-4o (OpenAI et al., 2024) and o3 (OpenAI, 2025b)) and ask them to predict and rank.

Verification task settings We construct the verification benchmark by selecting model–benchmark pairs from the artifact graph that include reported

Method	Link Prediction		Link Ranking	
	F1 \uparrow	Acc \uparrow	R@5 \uparrow	P@5 \uparrow
Metadata	25.9	49.0	47.3	21.1
MF	29.5	40.5	37.6	19.1
CN	62.7	81.0	81.0	39.1
AA	67.2	84.5	81.0	39.1
GPT-4o	59.1	79.1	89.2	37.2
GPT-4o + graph	61.9	81.7	91.7	38.6
GPT-o3	59.1	80.4	90.3	38.3
GPT-o3 + graph	64.2	83.8	92.6	40.2
NeoGNN	63.8	91.1	43.3	16.0
GATv2Conv	83.0	95.0	71.9	33.4

Table 1: **Link prediction and ranking results.** For link prediction, we sample 5 negative samples for each datapoint with preferential attachment that chooses the highest degree nodes. For link ranking, we construct a benchmark-centric ranking task by randomly sampling 10 negative model candidates per benchmark. More details in Section §6.

Method	Attr Prediction		Attr Ranking
	MAE \downarrow	RMSE \downarrow	NDCG@1 \uparrow
Global mean	0.159	0.201	–
Local mean	0.107	0.156	–
GPT-4o	0.136	0.216	0.468
GPT-4o + graph	0.077	0.171	0.459
GPT-o3	0.168	0.260	0.433
GPT-o3 + graph	0.060	0.126	0.431
GATv2Conv	0.071	0.128	0.519

Table 2: **Attribute prediction and ranking results.** For attribute prediction, we focus on the ‘accuracy’, ‘F1’, ‘BLEU’, ‘chrF’, ‘rouge’ metrics for calculation in the test split to make sure the prediction range is unified. For attribute ranking, we build each ranking task for each dataset under one type of metric. More details in Section §6.

scores (e.g., F1, accuracy, BLEU, ROUGE). For each pair, we test whether the reported result can be reproduced through execution-based evaluation, given the model and benchmark specifications. We collect a subset of 93 model-dataset pairs and avoid datasets that are too expensive to run and models that are too large to download.

Verification baseline settings As baselines, we adopt single-turn code generation and multi-turn agentic CodeAct (Wang et al., 2024). Importantly, this baseline does not include HuggingFace-specific design choices, but simply uses a coding sandbox to conduct experiments automatically. Our method is built based on smolagent and, for a fair comparison, we use GPT-5.2 (OpenAI, 2025a) as the underlying backbone model, set max turn number to be 10, timeout for each turn to 900s.

7 Prediction Task Evaluation Results

Graph structure consistently improves LLM prediction Augmenting LLMs with graph neighborhood information yields consistent gains across

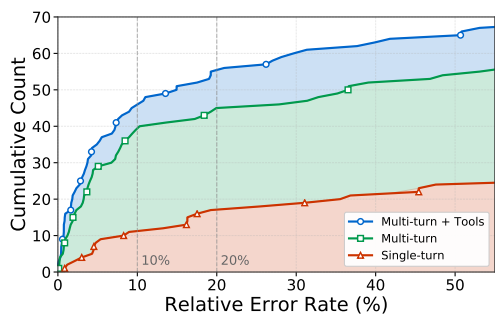


Figure 3: **Verification error rate distribution.** We calculate the cumulative count of datapoints with different error rates. Multi-turn and tool use mainly add more low-error-rate cases.

tasks (Tables 1 and 2), especially for link and attribute prediction. For example, adding 1-hop context improves GPT-o3’s link prediction F1 from 59.1 to 64.2 and reduces attribute MAE from 0.168 to 0.060. Ranking gains are smaller, likely due to strong semantic priors in LLMs. Under inductive splits with unseen models and datasets, performance degrades only mildly (e.g., GPT-4o+graph MAE increases from 0.07 to 0.08), indicating reasonable inductive robustness.

GNNs excel in transductive settings but generalize less inductively A lightweight 3-layer GATv2Conv with Voyage-3 embeddings achieves competitive or superior performance in the transductive setting as shown in Tables 1 and 2, matching GPT-4o+graph on attribute prediction and outperforming all baselines on link prediction. However, under inductive splits, GNN performance drops more substantially (MAE from 0.07 to 0.12), reflecting stronger dependence on observed graph structure compared to LLM-based methods.

Overall While current models are promising, significant headroom remains across all sub-tasks. These tasks go beyond metric evaluation and serve as proxies for scientific compatibility reasoning. Strong link and attribute prediction is a prerequisite for scalable automatic discovery, enabling agents to autonomously identify high-potential model–dataset combinations in the wild, highlighting the value of our task design.

8 Verification Task Evaluation Results

Multi-turn tool use enhances reproduction reliability As shown in Table 3, both multi-turn interactions and tool-use capabilities significantly improve the execution and reproduction rates of coding agents. Furthermore, Figure 3 illustrates that HuggingFace-specific tools effectively shift the distribution toward lower error rates, while the

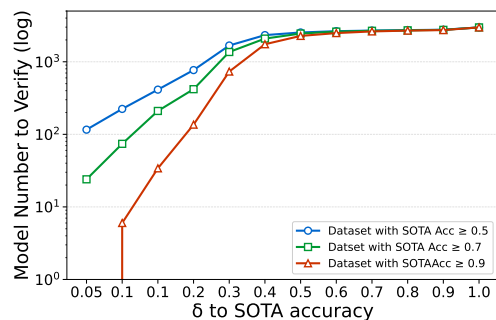


Figure 4: **Relationship between δ tolerance and the median number of candidate models to verify per dataset after prediction.** Given a dataset with known best accuracy, we use the GNN to predict accuracy for all 2,977 models and select candidates within $\delta\%$ of the existing SOTA scores.

frequency of high-error cases remains unchanged. Specifically, enabling agents to search and fetch HuggingFace model cards and metadata yields a clear performance boost (3.2 points in completion rate and 4.3 points in success rate). We attribute this to the agent’s ability to access model cards, which allows it to follow official code snippets and usage examples, thereby avoiding incorrect implementation. Despite these gains, several failure modes persist. The most common execution failures include missing result files (24%), zero-valued outputs (24%), and CUDA-related errors (12%). In addition, prompt sensitivity remains a key bottleneck: suboptimal evaluation prompts generated by our proposed coding agents can lead to unstable or degraded reproduced evaluation scores, particularly when evaluating LM candidates.

Overall Despite HuggingFace’s robust execution infrastructure, current coding agents reproduce reported results only 45.2% of the time, indicating limited reliability. This reveals verification on the artifact graph as both a critical benchmark of agent capability and a largely unsolved problem, highlighting robust auto-evaluation agents as a key direction for future research.

Metric	Single-turn	Multi-turn	MT + tools
Completed (%)	29.0	77.4	80.6
Succeeded (%)	11.8	40.9	45.2
#turn / dp	1	4.2	4.7
#tool-use / dp	–	–	10.4

Table 3: **Verification performance with coding agents.** “MT+tools” represents our proposed multi-turn tool-use agent. “/dp” refers to per datapoint. We report task completion rate and success rates. The success rate indicates that reproduced results are in 10% difference with the reported results.

9 Discussion

Q1: What signals drive effective prediction in the artifact graph? Effective prediction arises

MNLI	.94	.92	.91	.31	.92	.93	.84	.84	.30	.47	.89	.79	.56	.41	.50	.29	.60	.33	.27	.32	.30
SNLI	.88	.89	.93	.95	.92	.91	.71	.77	.95	.38	.85	.33	.60	.37	.38	.91	.40	.17	.37	.33	.34
ANLI	.67	.50	.64	.38	.56	.52	.61	.24	.82	.28	.30	.33	.17	.40	.38	.32	.42	.36	.26	.35	.35
XNLI	.30	.83	.87	.94	.88	.30	.62	.81	.36	.34	.28	.77	.07	.45	.37	.25	.30	.41	.30	.29	.30
RTE	.89	.18	.11	.86	.14	.17	.19	.40	.41	.84	.14	.22	.74	.56	.32	.18	.15	.54	.50	.50	.19
WNLI	.61	.45	.54	.49	.47	.54	.47	.43	.57	.47	.49	.45	.60	.43	.46	.52	.56	.57	.56	.43	.43
QNLI	.49	.76	.49	.49	.48	.95	.56	.50	.55	.67	.47	.50	.51	.61	.72	.52	.51	.50	.51	.51	.49
	DeBERTa-L-NLI	ModernBERT-B	ModernBERT-L	RoBERTa-L-ZS	DeBERTa-B-NLI	BART-L-MNLI	mDeBERTa-XNLI	DistilBERT	DeBERTa-L-ZS	Gemma-3-1B	RoBERTa-L-MNLI	BART-Yahoo	XtremeDistil	SGPT-125M	Gemma-3-1B-U	DistilBART	DeBERTa-B-ZS	Gemma-2B	Gemma-2-2B	GPT-Neo-1.3B	DeBERTa-B

Figure 5: **Verification results on NLI tasks.** We show the accuracy verification results conducted by the ARTIFACTLINKER with 21 models and 7 NLI datasets. The results are conducted with ARTIFACTLINKER and might have false negatives existed due to limited abilities of the coding agent.

from the interaction between node attributes and graph structure. LLM-based predictors use attributes for analogical reasoning, while GNNs encode them via Voyage-3 embeddings and exploit dense connectivity through message passing. GNNs perform best on high-degree nodes, LLMs excel on medium-degree nodes (approximately 15–30), and both struggle on low-degree nodes. Randomizing text embeddings significantly degrades attribute prediction RMSE. All predictions rely on LLM-summarized model cards with ground-truth labels removed to avoid leakage.

Q2: How does prediction reduce verification costs? We analyze how many models need to be verified per dataset after prediction-based filtering. For datasets with SOTA accuracy ≥ 0.7 , we count models predicted within a certain ratio of the ground-truth SOTA in Figure 4. At 5% tolerance, the median is 24 candidate models—a 99% reduction from the total. At 20% tolerance, the median rises to 419. This suggests 20-80 experiments typically suffice to identify near-optimal models.

10 End-to-end Case Study

To demonstrate the effectiveness of ARTIFACTLINKER when used to uncover real unseen links, we conduct a study on the classic task of Natural Language Inference (NLI). We consider seven representative benchmarks: MNLI (Williams et al., 2018), SNLI (Bowman et al., 2015), ANLI (Nie et al., 2020), XNLI (Conneau et al., 2018), RTE (Bentivogli et al., 2009), WNLI, and QNLI (Wang et al., 2018). ARTIFACTLINKER first predicts promising model–dataset pairs over the artifact graph, after which we select 21 relatively small models for fully automated coding-based verification. Results are shown in Figure 5.

SOTA discovery All evaluations in Figure 5 are conducted fully automatically without human in-

tervention. While some scores may be imperfect due to limitations of the coding agent, the results nonetheless reveal several meaningful discoveries. Notably, ARTIFACTLINKER identifies a new state-of-the-art result on SNLI: deberta-v3-large-zeroshot-v2.0¹, which does not report official SNLI results, achieves an accuracy of 0.95. We first verify that this score is consistent with the state-of-the-art estimates derived from our artifact graph. Furthermore, we confirm that it is similar to the best result reported on the SNLI leaderboard². We also include reproducible scripts in Appendix. E.

Research insight Beyond SOTA discovery, we observe strong performance correlation between MNLI and SNLI across multiple models, strong QNLI performance from BART-large-mnli³, and unexpectedly strong RTE results from Gemma-3-1B. Overall, this case study demonstrates ARTIFACTLINKER’s ability to efficiently discover SOTA-level results and uncover research insights not explicitly reported in HF, highlighting its potential as a scalable tool for automatic research.

11 Conclusion

We introduce ARTIFACTLINKER, a framework that enables automatic SOTA discovery on HuggingFace by modeling models and benchmarks as an artifact graph. Using a two-stage predict-and-verify approach—graph-based candidate prediction followed by execution-based verification. We further propose ARTIFACTBENCH to evaluate both stages, and show that ARTIFACTLINKER can discover high-performing links and reproduce many benchmark results automatically, pointing toward scalable and continual scientific discovery.

¹<https://huggingface.co/MoritzLaurer/deberta-v3-large-zeroshot-v2.0>

²<https://nlp.stanford.edu/projects/snli/>

³<https://huggingface.co/facebook/bart-large-mnli>

610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659

Limitations

Computational considerations in verification

While our two-stage framework effectively reduces the search space through candidate filtering in the prediction phase, the verification stage requires actual code execution for validation. For certain large-scale experiments or resource-intensive benchmarks, full reproduction may require non-trivial computational costs depending on available hardware resources. We note that our current implementation successfully handles the majority of common benchmarks, though scaling to extremely high-throughput scenarios across diverse hardware environments remains an interesting direction for future optimization.

Evaluation metric coverage Our framework currently emphasizes objective, quantitative metrics such as Accuracy, F1, and Exact Match, which are widely reported in model documentation and amenable to automatic extraction. These metrics cover a substantial portion of commonly used evaluation schemes in NLP. While our approach could potentially be extended to incorporate tasks involving human evaluation or qualitative assessment, we leave the integration of such subjective metrics as a natural extension for future work, as they would require additional methodological considerations for automated processing.

Ethics Consideration

Our work focuses on building an automatic discovery framework over open-source artifacts available on HuggingFace. All experiments are conducted exclusively on publicly released models and benchmarks that are freely accessible to the research community. We do not introduce any new human or sensitive data, nor do we attempt to deanonymize or misuse existing artifacts. The goal of our framework is to advance automated, reproducible, and scalable scientific discovery, and to help researchers more efficiently identify promising model–benchmark interactions. Nevertheless, we acknowledge that automated benchmarking may propagate existing biases and limitations present in the underlying models and datasets. To mitigate this, we emphasize transparency in data collection and reproducibility in our verification pipeline. We encourage the community to view our work as a step toward building more reliable and responsible auto-discovery systems, rather than a replacement for human oversight.

References

Adem Ait, Javier Luis Cánovas Izquierdo, and Jordi Cabot. 2023. [On the suitability of hugging face hub for empirical studies](#). *ArXiv*, abs/2307.14841.

Joeran Beel, Min-Yen Kan, and Moritz Baumgart. 2025. [Evaluating sakana’s ai scientist for autonomous research: Wishful thinking or an emerging reality towards ‘artificial research intelligence’ \(ari\)?](#) *ArXiv*, abs/2502.14297.

Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. 2009. The fifth pascal recognizing textual entailment challenge. *TAC*, 7(8):1.

Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. 2008. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment*, 2008(10):P10008.

Ben Bogin, Kejuan Yang, Shashank Gupta, Kyle Richardson, Erin Bransom, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2024. Super: Evaluating agents on setting up and executing tasks from research repositories. *Proceedings of EMNLP*.

Samuel Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 conference on empirical methods in natural language processing*, pages 632–642.

Shaked Brody, Uri Alon, and Eran Yahav. 2021. How attentive are graph attention networks? *arXiv preprint arXiv:2105.14491*.

Keith T Butler, Daniel W Davies, Hugh Cartwright, Olexandr Isayev, and Aron Walsh. 2018. Machine learning for molecular and materials science. *Nature*, 559(7715):547–555.

Joel Castaño, Rafael Cabañas, Antonio Salmeron, David Lo, and Silverio Martínez-Fernández. 2024. [How do machine learning models change?](#) *ArXiv*, abs/2411.09645.

Joel Castaño, Silverio Martínez-Fernández, Xavier Franch, and Justus Bogner. 2023. [Analyzing the evolution and maintenance of ml models on hugging face](#). *2024 IEEE/ACM 21st International Conference on Mining Software Repositories (MSR)*, pages 607–618.

Mark Chen. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.

Qiaosheng Chen, Kaijia Huang, Xiaofang Zhou, Weiqing Luo, Yuanning Cui, and Gong Cheng. 2025. [Benchmarking recommendation, classification, and tracing based on hugging face knowledge graph](#). In *Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*.

660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712

713	Junyan Cheng, Peter Clark, and Kyle Richardson. 2025. Language modeling by language models. <i>Proceedings of NeurIPS</i> .	767
714		768
715		769
716	Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. Xnli: Evaluating cross-lingual sentence representations. In <i>Proceedings of the 2018 conference on empirical methods in natural language processing</i> , pages 2475–2485.	770
717		771
718		772
719		773
720		774
721		
722	Ned Cooper, Tiffanie N. Horne, Gillian R. Hayes, Courtney Heldreth, Michal Lahav, Jess Holbrook, and Lauren Wilcox. 2022. A systematic review and thematic analysis of community-collaborative approaches to computing research. <i>Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems</i> .	775
723		776
724		777
725		778
726		779
727		
728		
729	Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In <i>2009 IEEE conference on computer vision and pattern recognition</i> , pages 248–255. Ieee.	780
730		781
731		782
732		783
733		
734	Sarah Gao and Andrew Gao. 2023. On the origin of llms: An evolutionary tree and graph for 15, 821 large language models. <i>ArXiv</i> , abs/2307.09793.	784
735		785
736		786
737	R. Heumüller, Sebastian Nielebock, J. Krüger, and F. Ortmeier. 2020. Publish or perish, but do not forget your software artifacts. <i>Empirical Software Engineering</i> , 25:4585 – 4616.	787
738		788
739		
740		
741	Peter Alexander Jansen, Marc-Alexandre Côté, Tushar Khot, Erin Bransom, Bhavana Dalvi, Bodhisattwa Prasad Majumder, Oyvind Tafjord, and Peter Clark. 2024. Discoveryworld: A virtual environment for developing and evaluating automated scientific discovery agents. <i>ArXiv</i> , abs/2406.06769.	789
742		790
743		791
744		792
745		793
746		
747	Peter Alexander Jansen, Oyvind Tafjord, Marissa Radensky, Pao Siangliulue, Tom Hope, Bhavana Dalvi, Bodhisattwa Prasad Majumder, D. S. Weld, and Peter Clark. 2025. Codescientist: End-to-end semi-automated scientific discovery with code-based experimentation. <i>ArXiv</i> , abs/2503.22708.	794
748		795
749		796
750		797
751		798
752		799
753	Seth Johnson, F. Samsel, G. Abram, Daniel L. Olson, Andrew J. Solis, Bridger Herman, P. Wolfram, C. Lenglet, and Daniel F. Keefe. 2019. Artifact-based rendering: Harnessing natural and traditional visual media for more expressive and engaging 3d visualizations. <i>IEEE Transactions on Visualization and Computer Graphics</i> , 26:492–502.	800
754		801
755		802
756		803
757		804
758		805
759		806
760	Gyeongwon James Kim, Alex Wilf, Louis philippe Morency, and Daniel Fried. 2025. From reproduction to replication: Evaluating research agents with progressive code masking. <i>ArXiv</i> , abs/2506.19724.	807
761		808
762		
763		
764	H. Kitano. 2021. Nobel turing challenge: creating the engine for scientific discovery. <i>NPJ Systems Biology and Applications</i> , 7.	809
765		810
766		
	Stefan Kramer, Mattia Cerrato, S. Džeroski, and R. King. 2023. Automated scientific discovery: From equation discovery to autonomous discovery systems. <i>ArXiv</i> , abs/2305.02251.	811
		812
		813
	Benjamin Laufer, Hamidah Oderinwale, and Jon Kleinberg. 2025. Anatomy of a machine learning ecosystem: 2 million models on hugging face. <i>ArXiv</i> , abs/2508.06811.	814
		815
		816
		817
	C. J. Lissa, A. Brandmaier, Loek Brinkman, Anna-Lena Lamprecht, Aaron Peikert, Marijn E. Struiksmā, and Barbara M. I. Vreede. 2020. Wocrs: A workflow for open reproducible code in science. <i>Data Sci.</i> , 4:29–49.	818
		819
		820
		821
	Chris Lu, Cong Lu, R. T. Lange, J. Foerster, Jeff Clune, and David Ha. 2024. The ai scientist: Towards fully automated open-ended scientific discovery. <i>ArXiv</i> , abs/2408.06292.	
	T. Marić, Dennis Gläser, Jan-Patrick Lehr, Ioannis Pagiannidis, B. Lambie, Christian H. Bischof, and Dieter Bothe. 2023. A pragmatic workflow for research software engineering in computational science. <i>ArXiv</i> , abs/2310.00960.	
	Yanzhou Mu, Rong Wang, Juan Zhai, Chunrong Fang, Xiang Chen, Jiacong Wu, An Guo, Jiawei Shen, Bingzhuo Li, and Zhenyu Chen. 2025. Understanding llm-centric challenges for deep learning frameworks: An empirical analysis. In <i>unknown</i> .	
	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial nli: A new benchmark for natural language understanding. In <i>Proceedings of the 58th annual meeting of the association for computational linguistics</i> , pages 4885–4901.	
	OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, and 401 others. 2024. Gpt-4o system card. <i>Preprint</i> , arXiv:2410.21276.	
	OpenAI. 2025a. Introducing OpenAI gpt-5.2. Accessed: 2026-01-06.	
	OpenAI. 2025b. Introducing OpenAI o3 and o4-mini. Accessed: 2026-01-06.	
	Mohammad Shahedur Rahman, Peng Gao, and Yuede Ji. 2025. Hugginggraph: Understanding the supply chain of llm ecosystem. <i>ArXiv</i> , abs/2507.14240.	
	Marwin HS Segler, Mike Preuss, and Mark P Waller. 2018. Planning chemical syntheses with deep neural networks and symbolic ai. <i>Nature</i> , 555(7698):604–610.	
	Minju Seo, Jinheon Baek, Seongyun Lee, and Sung Ju Hwang. 2025. Paper2code: Automating code generation from scientific papers in machine learning. <i>ArXiv</i> , abs/2504.17192.	

822	D. Serrano, F. C. Luciano, B. J. Anaya, Baris On-	Adina Williams, Nikita Nangia, and Samuel Bowman.	878
823	goren, Aytug Kara, Gracia Molina, Bianca I Ramirez,	2018. A broad-coverage challenge corpus for sen-	879
824	Sergio A Sánchez-Guiraes, Jesus A. Simon, Greta	tence understanding through inference. In <i>Proceed-</i>	880
825	Tomietto, Chrysi Rapti, Helga K. Ruiz, Satyavati	<i>ings of the 2018 conference of the North American</i>	881
826	Rawat, Dinesh Kumar, and A. Lalatsa. 2024. Artificial intelligence (ai) applications in drug discovery and drug delivery: Revolutionizing personalized medicine. <i>Pharmaceutics</i> , 16.	<i>chapter of the association for computational linguistics: human language technologies, volume 1 (long papers)</i> , pages 1112–1122.	882
827			883
828			884
829			
830	Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang,	Yanzheng Xiang, Hanqi Yan, Shuyin Ouyang, Lin Gui,	885
831	Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng Li,	and Yulan He. 2025. Scireplicate-bench: Benchmarking llms in agent-driven algorithmic reproduction from research papers. <i>ArXiv</i> , abs/2504.00255.	886
832	and Y. Zhuang. 2023. Taskbench: Benchmarking large language models for task automation. <i>ArXiv</i> ,		887
833	abs/2311.18760.		888
834			
835	Zachary S. Siegel, Sayash Kapoor, Nitya Nagdir,	Tian Xie and Jeffrey C Grossman. 2018. Crystal graph	889
836	Benedikt Stroebel, and Arvind Narayanan. 2024.	convolutional neural networks for an accurate and	890
837	Core-bench: Fostering the credibility of published research through a computational reproducibility agent benchmark. <i>Trans. Mach. Learn. Res.</i> , 2024.	interpretable prediction of material properties. <i>Physi-</i>	891
838		<i>cal review letters</i> , 120(14):145301.	892
839			
840	Kanishka Silva, Marcel R. Ackermann, Heike Fliegl,	Xinyu Yang, Weixin Liang, and James Zou. 2024. Navigating dataset documentations in ai: A large-scale analysis of dataset cards on hugging face. <i>ArXiv</i> ,	893
841	G. Gesese, Fidan Limani, Philipp Mayr, Peter	abs/2401.13822.	894
842	Mutschke, A. Oelen, Muhammad Asif Suryani,		895
843	Sharmila Upadhyaya, Benjamin Zopilko, Harald	Yujie You, Xin Lai, Ying Pan, Huiru Zheng, J. Vera,	897
844	Sack, and Stefan Dietze. 2025. Research knowledge graphs in nfidi4datascience: Key activities, achievements, and future directions. <i>ArXiv</i> , abs/2508.02300.	Suran Liu, Senyi Deng, and Le Zhang. 2022. Artificial intelligence in cancer target identification and drug discovery. <i>Signal Transduction and Targeted</i>	898
845		<i>Therapy</i> , 7.	899
846			900
847	Giulio Starace, Oliver Jaffe, Dane Sherburn, James	Haofei Yu, Zhaochen Hong, Zirui Cheng, Kunlun	902
848	Aung, Jun Shern Chan, Leon Maksin, Rachel Dias,	Zhu, Keyang Xuan, Jinwei Yao, Tao Feng, and	903
849	Evan Mays, Benjamin Kinsella, Wyatt Thompson,	Jiaxuan You. 2024. Researchtown: Simulator of human research community. <i>arXiv preprint</i>	904
850	and 1 others. 2025. Paperbench: Evaluating ai’s ability to replicate ai research. <i>arXiv preprint</i>	<i>arXiv:2412.17767</i> .	905
851	<i>arXiv:2504.01848</i> .		906
852			
853	Jonathan M Stokes, Kevin Yang, Kyle Swanson, Wen-	Seongjun Yun, Seoyoon Kim, Junhyun Lee, Jaewoo	907
854	gong Jin, Andres Cubillos-Ruiz, Nina M Donghia,	Kang, and Hyunwoo J Kim. 2021. Neo-gnns: Neigh-	908
855	Craig R MacNair, Shawn French, Lindsey A Carfrae,	borhood overlap-aware graph neural networks for	909
856	Zohar Bloom-Ackermann, and 1 others. 2020. A	link prediction. <i>Advances in Neural Information Pro-</i>	910
857	deep learning approach to antibiotic discovery. <i>Cell</i> ,	<i>cessing Systems</i> , 34:13683–13694.	911
858	180(4):688–702.		
859	R. Urbanowicz, Robert F. Zhang, Yuhan Cui, and Pran-		
860	shu Suri. 2022. Streamline: A simple, transparent, end-to-end automated machine learning pipeline facilitating data analysis and algorithm comparison. In <i>Genetic Programming Theory and Practice</i> .		
861			
862			
863			
864	A. Vişan and I. Neğüt. 2024. Integrating artificial intelligence for drug discovery in the context of revolutionizing drug delivery. <i>Life</i> , 14.		
865			
866			
867	Alex Wang, Amanpreet Singh, Julian Michael, Felix		
868	Hill, Omer Levy, and Samuel Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In <i>Proceedings of the 2018 EMNLP workshop BlackboxNLP: Analyzing and interpreting neural networks for NLP</i> , pages 353–		
869			
870			
871			
872			
873			
874	Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang,		
875	Yunzhu Li, Hao Peng, and Heng Ji. 2024. Executable code actions elicit better llm agents. In <i>Forty-first International Conference on Machine Learning</i> .		
876			
877			

912	A Potential Risks		
913	Metadata Quality Considerations.	Our prediction	961
914	stage uses metadata from the HuggingFace	community, including model descriptions, architec-	962
915	tural specifications, and training configurations. As	with any crowdsourced platform, there may occa-	963
916	sionally be instances of incomplete or imprecise	documentation that could affect initial ranking pre-	
917	dictions.		
918	Performance-Oriented Discovery Scope.	ARTIFACTLINKER is designed to identify high-	
919	performing model-benchmark combinations based	on standard evaluation metrics. As with any auto-	
920	mated performance discovery tool, users should	exercise standard research practices by conduct-	
921	ing appropriate safety and ethical assessments be-	fore deploying discovered models in production	
922	environments, particularly for applications involv-	ing user-facing content generation. We view our	
923	framework as a research tool that augments human	decision-making rather than replacing it. Users	
924	retain full responsibility for evaluating whether dis-	covered models meet the safety, fairness, and eth-	
925	ical requirements of their specific use cases, and	we encourage comprehensive evaluation beyond	
926	performance metrics alone.		
927			
928			
929			
930			
931			
932			
933			
934			
935			
936			
937			
938	B Scientific Artifacts		
939	B.1 Data license		
940	All data in ARTIFACTBENCH are derived from the	HuggingFace Hub, which hosts open-source mod-	
941	els and benchmarks under various public licenses	(<i>e.g.</i> , Apache 2.0, MIT, and CC-BY). We plan to re-	
942	lease ARTIFACTBENCH under the Open Database	License (ODbL), which permits use, sharing, and	
943	modification of the data while requiring proper at-	tribution and that any derivative works remain open-	
944	ly available under the same license.		
945			
946			
947			
948			
949	B.2 Model license		
950	Our work relies on multiple foundation models, in-	cluding chatgpt-4o-latest and o3-2025-04-16.	
951	Specifically, we access chatgpt-4o-latest and	texttto3-2025-04-16 via the official OpenAI API.	
952	We utilize the official inference API provided by	VoyageAI to use voyage-3.	
953	The chatgpt-4o-latest, o3-2025-04-16,	and voyage-3 models are closed-source and	
954	operate under proprietary licenses. We use them	only for academic, and non-commercial purposes	
955	and ensure all inputs come from publicly available	data, complying with their usage restrictions. We	961
956		make no modifications to these models and use	962
957		them as-is via their public APIs.	963
958			
959			
960			
	B.3 Data Usage		964
	All data used in this study, including model speci-	fications, benchmark descriptions, and evaluation	965
	results, were collected from the publicly available	HuggingFace Hub. Our data collection focuses ex-	966
	clusively on scientific artifacts (models, datasets,	and benchmarks) and does not include any person-	967
	ally identifiable information. All collected data	are derived from publicly released resources that	968
	are freely accessible to the research community,	in compliance with the HuggingFace platform’s	969
	Terms of Service.		970
			971
			972
			973
			974
			975
	C Experimental Details		976
	C.1 GNN-based Modeling and Training		977
	Details		978
	In this section, we describe the architecture of our	GNN-based model and provide details on its train-	979
	ing setup for both link prediction and attribute re-	gression tasks.	980
			981
			982
	C.2 Model Architecture		983
	Our link prediction model, GNNLinkPredictor,	combines a graph neural network encoder with	984
	a flexible edge-scoring decoder. The encoder is	based on GATv2 and stacks multiple convolutional	985
	layers with GraphNorm, PReLU activations, and fea-	ture dropout. Residual connections are added when	986
	input and output dimensions match, and self-loops	as well as edge dropout can be applied for sta-	987
	bility. To integrate information across layers, we	adopt Jumping Knowledge (JK), supporting con-	988
	catenation with projection, element-wise max pool-	ing, or using the last layer only. The resulting	989
	node embeddings are passed into the edge decoder	(EdgePredictor), which supports various scoring	990
	functions, including dot product, cosine similar-	ity, bilinear transformation, linear projection on	991
	concatenated features, and a two-layer MLP. This	design balances expressiveness and flexibility: the	992
	encoder captures rich multi-hop structure, while	the decoder adapts to diverse relational patterns.	993
			994
			995
			996
			997
			998
			999
			1000
			1001
			1002
			1003
	C.3 Training for Link Prediction		1004
	For binary edge classification, we train our GNN	model over positive and negative edges sampled	1005
	from the artifact graph, using pre-computed node		1006
			1007

embeddings as input. We optimize with AdamW (learning rate 5×10^{-3} , weight decay 10^{-4}) and use a ReduceLROnPlateau scheduler (patience 10, decay factor 0.8, minimum learning rate 10^{-6}). Training runs for up to 300 epochs with early stopping (patience 40) based on validation AUC, and class imbalance is corrected by weighting the binary cross-entropy loss with $\frac{N_{\text{neg}}}{N_{\text{pos}}}$. Mixed-precision training is enabled on GPU. We evaluate every 50 epochs, select the best checkpoint by validation AUC, and report accuracy, precision, recall, F1, ROC-AUC, and average precision on the test set. Unless otherwise noted, we use 64 hidden dimensions, 3 layers, 3 heads, dropout 0.2, and seed 42.

C.4 Training for Attribute Regression

For continuous edge-attribute prediction, we train on positive edges with metric values, normalizing values greater than 1 to $[0, 1]$. Node embeddings are pre-computed and message passing uses the training graph’s edges. The model regresses edge scores directly in logit space: ground-truth targets are transformed via $\text{logit}(y) = \log \frac{y}{1-y}$, and MSE is computed against the raw decoder logits. This stabilizes training near the boundaries. We use Adam (learning rate 0.005, weight decay 10^{-5}), Xavier-uniform initialization for linear layers, dropout 0.2, hidden size 128, 3 layers, and 8 heads (GATv2 backbone). Training runs up to 500 epochs, with frequent evaluation early (every 10 epochs for the first 50, then every 25), and the best checkpoint is selected by lowest validation MSE. At evaluation, logits are clipped to $[-10, 10]$ and mapped through sigmoid to report MSE, MAE, RMSE, R^2 , and MAPE.

C.5 Model Size and Budget

For the close-sourced models, the model size is unknown. The GNN model we used is within 100 MB. The budget for conducting prediction/ranking/verification is $< \$1000$.

D The Use of Large Language Models (LLMs)

We used ChatGPT as a writing assistant to help us write part of the paper. Additionally, we utilize the power of CodePilot to help us code faster. However, all the AI-generated writing and coding components are manually checked and modified. There is no full AI-generated content in the paper.

E SOTA Discovery Details

We provide the discovered state-of-the-art evaluation code here:

Listing 1: Automatic SNLI evaluation for SOTA verification.

```

import json
import random
import torch
from datasets import load_dataset
from transformers import (
    AutoTokenizer,
    AutoConfig,
    AutoModelForSequenceClassification
)
MODEL_ID =
    "MoritzLaurer/deberta-v3-large-zeroshot
-v2.0"
DATASET_ID = "stanfordnlp/snli"
SPLIT = "validation"
MAX_SAMPLES = 20000
BATCH_SIZE = 16
MAX_LENGTH = 256
SEED = 42

random.seed(SEED)
torch.manual_seed(SEED)

dataset = load_dataset(DATASET_ID,
    split=SPLIT)

def is_valid(example):
    return (
        example["label"] is not None
        and example["label"] != -1
        and example["premise"]
        and example["hypothesis"]
    )

valid_indices = [i for i in
    range(len(dataset)) if
    is_valid(dataset[i])]
random.shuffle(valid_indices)
dataset =
    dataset.select(valid_indices[:MAX_SAMPLE
S])

tokenizer =
    AutoTokenizer.from_pretrained(MODEL_ID,
    use_fast=True)
config =
    AutoConfig.from_pretrained(MODEL_ID)
model = AutoModelForSequenceClassification.
    from_pretrained(MODEL_ID)

device = torch.device("cuda" if
    torch.cuda.is_available() else "cpu")
model.to(device)
model.eval()

# Determine model label mapping
if getattr(config, "label2id", None) and
    len(config.label2id) > 0:
    label2id = {k.lower(): int(v) for k, v
    in config.label2id.items()}
else:
    label2id = {v.lower(): int(k) for k, v

```

```

1121         in config.id2label.items())
1122
1123     entailment_id = label2id["entailment"]
1124     not_entailment_id =
1125         label2id["not_entailment"]
1126
1127     snli_labels =
1128         dataset.features["label"].names
1129     snli_entailment_id =
1130         snli_labels.index("entailment")
1131
1132     correct, total = 0, 0
1133
1134     for i in range(0, len(dataset), BATCH_SIZE):
1135         batch = dataset[i:i+BATCH_SIZE]
1136
1137         enc = tokenizer(
1138             batch["premise"],
1139             batch["hypothesis"],
1140             padding=True,
1141             truncation=True,
1142             max_length=MAX_LENGTH,
1143             return_tensors="pt"
1144         )
1145         enc = {k: v.to(device) for k, v in
1146             enc.items()}
1147
1148         labels = [
1149             entailment_id if y ==
1150             snli_entailment_id else
1151             not_entailment_id
1152             for y in batch["label"]
1153         ]
1154
1155         with torch.no_grad():
1156             preds =
1157             model(**enc).logits.argmax(dim=-1).cpu().
1158             tolist()
1159
1160         for p, y in zip(preds, labels):
1161             correct += int(p == y)
1162             total += 1
1163
1164     accuracy = correct / total
1165     print("N_=", total)
1166     print("Accuracy_=", accuracy)
1167
1168     with open("results.json", "w") as f:
1169         json.dump({"accuracy": accuracy}, f)

```

F LLM-based Modeling Details

We directly rely on in-context learning for LLM-based modeling without training. We provide examples of our prompts for different tasks below:

Listing 2: LLM with graph link prediction example

```

1175     Given a machine learning model named
1176         'tensorblock/bloomz-3b-GGUF' and a
1177         dataset named 'SetFit/rte'.
1178
1179
1180     More information about this model: The
1181         model 'bloomz-3b1' is a multilingual
1182         text-generation model based on the
1183         BLOOM architecture. It is fine-tuned on
1184         a variety of datasets such as the
1185         BigScience xP3, which contributes to

```

```

its ability to handle multiple
languages and tasks. The model supports
a wide range of languages including
English, Spanish, French, Chinese, and
many others, making it highly versatile
for global applications. It is also
capable of understanding and generating
programming code in several languages,
like C++, Java, Python, and more. This
flexibility is reflected in its usage
across diverse tasks, from sentiment
analysis and question answering to
natural language inference and program
synthesis. The model operates under the
bigscience-bloom-rail-1.0 license,
ensuring its use in accordance with
open research and collaboration
principles.

```

More information about this dataset: The Glue RTE dataset is derived from the Recognizing Textual Entailment (RTE) task, a subset of the GLUE benchmark. It is designed to evaluate models on the task of natural language inference, which involves determining if one sentence logically follows from another. In the ported version hosted on HuggingFace's platform, it retains its original purpose but introduces some modifications for easier processing. Notably, the columns originally named 'sentence1' and 'sentence2' have been renamed to 'text1' and 'text2' respectively. This dataset is often used to train and evaluate machine learning models on their ability to understand and process natural language, particularly in contexts where discerning relationships between two text statements is crucial. One important characteristic of the dataset is that its test split is unlabeled, with all entries in the label column set to -1. This is in line with many benchmark datasets where the test labels are withheld to prevent overfitting and encourage robust model evaluation.

There are other models that are evaluated on the dataset to judge whether the model and dataset are connected:

- 42dot/42dot_LLM-PLM-1.3B: 42dot LLM-PLM is a pre-trained language model developed by 42dot, designed to handle both Korean and English text. It is part of the 42dot LLM series, leveraging a 1.3 billion parameter configuration based on a Transformer decoder architecture, akin to LLaMA 2. The model is trained with a corpus of diverse text sources, both in Korean and English, aimed at providing a versatile foundation for various natural language tasks. The architecture consists of 24 layers, 32 attention heads, a hidden size of 2048, and an FFN size of 5632. Pre-training

1396
1397
1398
1399
1400

evaluated on this dataset. Provide your answer as a JSON object with two keys: 'prediction' (a boolean, true or false) and 'reason' (a brief explanation of your reasoning).

Listing 3: LLM with graph attribute prediction example

1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463

Given a machine learning model named 'tasksource/ModernBERT-large-nli' and a dataset named 'stanfordnlp/snli'.

More information about this model:
ModernBERT is a multi-task fine-tuned model specifically designed for natural language inference (NLI) tasks. It has been trained on a diverse set of NLI datasets including MNLI, ANLI, SICK, WANLI, doc-nli, LingNLI, FOLIO, FOL-NLI, LogicNLI, and Label-NLI among others. The training was conducted over 200,000 steps using an Nvidia A30 GPU, resulting in a powerful model for reasoning tasks. ModernBERT surpasses the performance of models like llama 3.1 8B Instruct in specific tasks such as ANLI and FOLIO. It excels in long context reasoning, sentiment analysis, and zero-shot classification with new labels. The model is versatile, offering significant potential for fine-tuning on single tasks like SST, but it is particularly effective out-of-the-box for zero-shot classification and NLI tasks.

More information about this dataset: The Stanford Natural Language Inference (SNLI) corpus is a comprehensive dataset designed to support the task of natural language inference (NLI), also known as recognizing textual entailment (RTE). Version 1.0 of the dataset comprises 570,000 human-written sentence pairs labeled with entailment, contradiction, and neutral. These labels are crucial for training and evaluating models that aim to determine the inference relationship between two short texts. The SNLI dataset is an essential resource for developing systems that can understand and predict semantic relationships in natural language, making it a benchmark for text representation learning methodologies. The dataset was created through crowdsourcing efforts, with premises originating from the Flickr 30k and VisualGenome corpora, and hypotheses generated by crowdworkers on Amazon Mechanical Turk. Each sentence pair includes a premise and a hypothesis, and the task is to classify the relationship between them. The dataset provides a balanced classification challenge, offering a rich source of data for training machine learning models in natural language processing applications.

1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533

tasksource/ModernBERT-large-nli's performance on other datasets:
- pietrolesci/nli_fever: accuracy: 0.780 (info: The dataset in question is a modification of the FEVER dataset, tailored for Natural Language Inference (NLI) research. Originally, the FEVER dataset consisted of claims derived from Wikipedia, each paired with a label indicating the veracity of the claim. However, this setup did not align with standard NLI tasks which typically involve a pair of sequences (e.g., premise and hypothesis) mapped to a label. To bridge this gap and facilitate NLI research, the creators of this dataset have reformatted it to pair claims with textual evidence, thus converting it into a pair-of-sequences to label dataset. This transformation enables the application of NLI models on the FEVER dataset. The labels are mapped using predefined categories, where 'SUPPORTS' is mapped to entailment, 'NOT ENOUGH INFO' to neutral, and 'REFUTES' to contradiction. Additionally, the dataset includes a 'verifiable' column encoded to indicate whether a claim can be verified. Despite these modifications, the dataset maintains consistency with the original FEVER dataset, ensuring reliability and validity for research purposes.)
- yale-nlp/FOLIO: accuracy: 0.710 (info: The dataset is designed to provide a comprehensive collection of data relevant to the training and evaluation of machine learning models across various applications. It comprises a wide range of data types and structures, ensuring versatility and scalability for different research and development purposes. The dataset is curated to support both supervised and unsupervised learning tasks, facilitating experimentation with classification, regression, clustering, and more advanced machine learning methodologies. Each data entry in the dataset is meticulously labeled and documented, providing clear context and metadata to aid in the effective utilization of the data. The dataset is continually updated to incorporate the latest advancements and feedback from the user community, ensuring it remains relevant and useful. It is distributed under the MIT license, allowing for broad usage, adaptation, and distribution, making it an ideal resource for both academic and commercial projects. The dataset's accessibility and comprehensiveness make it a valuable asset for data scientists, researchers, and developers aiming to advance machine learning capabilities.)

1534	- rediska0123/puzzte: accuracy: 0.590	designed to evaluate the ability of	1604
1535	(info: The dataset is designed to serve	transformer models to perform reasoning	1605
1536	as a resource for binary	over natural language. This dataset	1606
1537	question-answering models. It consists	consists of synthetic data generated to	1607
1538	of questions that require a true or	simulate logical reasoning tasks using	1608
1539	false answer. The main features of the	controlled language. The primary	1609
1540	dataset include a 'question' field,	features of the dataset include	1610
1541	which contains the text of the question	'context', 'question', 'label', and	1611
1542	itself, and an 'answer' field, which is	'config', each represented as a string,	1612
1543	a boolean indicating whether the answer	which are structured to test the	1613
1544	is true or false. The dataset is	model's capacity to infer and deduce	1614
1545	available in a single configuration	truths from given premises. The dataset	1615
1546	named 'default' and contains a single	is provided in three splits: 'train'	1616
1547	split, 'train', which includes 911	with 480,152 examples, 'dev' with	1617
1548	examples. The total dataset size is	75,872 examples, and 'test' with	1618
1549	242,839 bytes, while the download size	151,911 examples, collectively	1619
1550	is 64,267 bytes.)	requiring a dataset size of 372,450,135	1620
1551	- causal-nlp/CLadder: accuracy: 0.890	bytes. It is licensed under the	1621
1552	(info: The dataset consists of multiple	Apache-2.0 license and primarily	1622
1553	configuration files specifying	intended for English language	1623
1554	different data splits. In the provided	processing. The dataset is based on the	1624
1555	configurations, there are two main data	work presented in the paper	1625
1556	files: 'full_v1.5_default' and	'Transformers as Soft Reasoners over	1626
1557	'full_v1'. These files are likely CSV	Language' by Peter Clark, Oyvind	1627
1558	formatted and are stored in the 'data'	Tafjord, and Kyle Richardson, which was	1628
1559	directory. The 'default' configuration	presented at the Twenty-Ninth	1629
1560	suggests that it might be the primary	International Joint Conference on	1630
1561	or recommended setup for working with	Artificial Intelligence in 2020.)	1631
1562	the dataset. Each configuration is	- tasksource/babi_nli: accuracy: 0.980	1632
1563	defined under the 'configs' key, with	(info: The babi_nli dataset is designed	1633
1564	each containing a 'config_name' and	for natural language inference tasks,	1634
1565	associated 'data_files'. The data files	specifically focusing on logical	1635
1566	are specified with a 'split' name and	reasoning capabilities. This dataset	1636
1567	the corresponding file path. This setup	supports various logic-based tasks by	1637
1568	allows users to select different	providing pairs of premises and	1638
1569	versions or subsets of the data	hypotheses, each labeled as either	1639
1570	depending on their use case,	'entailed' or 'not-entailed.' The	1640
1571	facilitating experiments with varied	dataset is monolingual, using English	1641
1572	data distributions or updated datasets.	language text, and is created through a	1642
1573	Such configuration flexibility is	combination of expert-generated	1643
1574	crucial for testing machine learning	annotations and crowdsourcing efforts.	1644
1575	models under different conditions or	It falls under the size category with	1645
1576	for ensuring reproducibility in	examples ranging between 1,000 and	1646
1577	research.)	10,000. The dataset offers multiple	1647
1578	- MoritzLaurer/dataset_train_nli: accuracy:	configurations, each targeting	1648
1579	0.950 (info: The dataset comprises	different aspects of logical reasoning	1649
1580	multiple features including 'text',	such as agents' motivations,	1650
1581	'hypothesis', 'labels', 'task_name',	coreference, deduction, induction, and	1651
1582	and 'label_text'. The 'labels' feature	several others. Each configuration	1652
1583	is a class label with two possible	provides its own train, validation, and	1653
1584	values: 'entailment' (0) and	test splits, with varying dataset sizes	1654
1585	'not_entailment' (1). This dataset is	and download sizes. The dataset can be	1655
1586	primarily used for natural language	used to evaluate models on their	1656
1587	inference tasks, where the goal is to	ability to perform natural language	1657
1588	determine if a given hypothesis is	inference tasks, an essential component	1658
1589	entailed by the premise text. The data	of many modern applications in natural	1659
1590	is organized into a single	language processing.)	1660
1591	configuration named 'default' and	stanfordnlp/snli's performance with other	1661
1592	focuses on a training split. The	models:	1662
1593	training data consists of 1,018,733	- tasksource/ModernBERT-base-nli: accuracy:	1663
1594	examples, occupying 315,017,473 bytes.	0.830 (info: ModernBERT is a multi-task	1664
1595	The total download size for the dataset	fine-tuned transformer model	1665
1596	is approximately 206 MB. This	specifically designed for natural	1666
1597	structured setup allows for easy	language inference (NLI) tasks. It	1667
1598	implementation of machine learning	incorporates datasets such as MNLI,	1668
1599	models aimed at understanding the	ANLI, SICK, WANLI, doc-nli, LingNLI,	1669
1600	entailment relationship between pairs	FOLIO, FOL-NLI, LogicNLI, and	1670
1601	of sentences.)	Label-NLI, among others. The model is	1671
1602	- tasksource/rulemaker: accuracy: 0.990	an 'instruct' version, optimized for	1672
1603	(info: The 'rulemaker' dataset is	tasks requiring reasoning and zero-shot	1673

1674 classification. ModernBERT's training
1675 regime involved 200,000 steps on an
1676 Nvidia A30 GPU, enhancing its
1677 capabilities in long-context reasoning
1678 and sentiment analysis. It surpasses
1679 other models like llama 3.1 8B Instruct
1680 in specific NLI tasks such as ANLI and
1681 FOLIO. Though additional performance
1682 improvements can be achieved through
1683 single-task fine-tuning (e.g., SST),
1684 the current version excels in zero-shot
1685 classification and natural language
1686 inference across various categories
1687 including contradiction, entailment,
1688 and neutral classification. The model's
1689 architecture employs different
1690 classification heads on top of the same
1691 transformer, making it versatile for
1692 multiple applications. Notably,
1693 ModernBERT benefits from NLI training
1694 data such as the 'label-nli' dataset,
1695 specifically designed to enhance
1696 zero-shot classification abilities.)

1697
1698 Please predict the accuracy that
1699 `tasksource/ModernBERT-large-nli` would
1700 achieve on `stanfordnlp/snli`. Provide
1701 your answer as a JSON object with two
1702 keys: 'prediction' (a float between 0
1703 and 1) and 'reason' (a brief
1704 explanation of your reasoning).