

Granite Foundation Models

IBM Research

Abstract—We introduce the Granite series of decoder-only foundation models for generative artificial intelligence (AI) tasks that are ready for enterprise use. We report on the architecture, capabilities, underlying data and data governance, training algorithms, compute infrastructure, energy and carbon footprint, testing and evaluation, socio-technical harms and mitigations, and usage policies.

Index Terms—foundation model, large language model, generative AI, data governance, contrastive fine-tuning, energy consumption, evaluation, socio-technical harms, usage governance, transparent documentation

I. INTRODUCTION

In this technical report, we present the Granite series of decoder-only foundation models for generative artificial intelligence (AI) tasks. The first in this series, granite.13b, is an English-only large language model (LLM). Using self-supervised learning, this base model has been trained on an IBM-curated pre-training dataset described in Section II. IBM relies on its internal end-to-end data and AI model lifecycle governance process and capabilities to develop enterprise-grade foundation models and is making similar capabilities available to customers of its watsonx platform.

The base model is the jumping-off point for two variants: granite.13b.instruct and granite.13b.chat. The first variant, granite.13b.instruct, has undergone supervised fine-tuning to enable better instruction following [1] so that the model can be used to complete enterprise tasks via prompt engineering. The second variant, granite.13b.chat, has undergone a novel contrastive fine-tuning after supervised fine-tuning to further improve the model’s instruction following, mitigate certain notions of harms, and encourage its outputs to follow certain social norms and have some notion of helpfulness [2]–[4]. We emphasize that these notions are not universal and discuss this point to a greater extent in Section VI on socio-technical harms and risks.

The granite.13b.instruct and granite.13b.chat models are made available by IBM through the watsonx platform [5]. IBM indemnifies customer use of these models on the watsonx platform, providing the same contractual intellectual property protections for IBM-developed AI models as it does for all of IBM’s products according to IBM Standard Terms and Conditions.

A. Overview of Capabilities

The 13b in the name indicates the model has 13 billion parameters. Furthermore, the base granite.13b decoder-only model has multi-query attention with learned position embeddings,

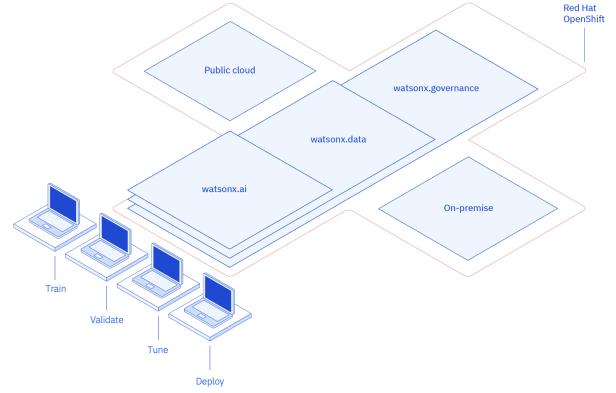


Fig. 1. A conceptual diagram of the watsonx platform.

has been trained on 1 trillion tokens created with the GPT-NeoX 20B tokenizer [6], and has a context length of 8 thousand tokens. As discussed in Section V, the Granite models are competitive in their ‘weight class’ on benchmark evaluations while being enterprise-ready in governance dimensions.

Some of the key enterprise tasks (common across sectors) for which the Granite models may be used are: retrieval-augmented generation, summarization, content generation, named entity recognition, insight extraction, and classification. The Granite models may be adapted to the specific tasks arising in particular enterprise applications through prompt engineering in the watsonx platform, which is illustrated in Fig. 1. Other series of models that IBM is developing are Sandstone: encoder-decoder models designed to be tuned for specific tasks and Obsidian: modular universal transformer models suitable for high inference efficiency.

B. Overview of the Granite Pre-Training Dataset

To support the training of large enterprise-grade foundation models, including granite.13b, IBM curated a massive dataset of relevant unstructured language data from sources across academia, the internet, enterprise (e.g., financial, legal), and code. In a rare move from a major provider of proprietary LLMs, IBM demonstrates its commitment to transparency and responsible AI by publishing descriptions of its training dataset in Section II.

The Granite pre-training dataset was created as a proprietary alternative to commonly used open-source data compilations for LLM training such as “The Pile” [7] or “C4” [8]. Some domains that are key for enterprise natural language processing are relatively under-represented in these compilations. Additionally these data compilations have been criticized for

containing toxic, harmful, or pirated content [9]. By curating our own pre-training data corpus, IBM takes significant steps towards addressing these and other issues.

The IBM curated pre-training dataset is continually growing and evolving, with additional data reviewed and considered to be added to the corpus at regular intervals. In addition to increasing the size and scope of pre-training data, new versions of these datasets are regularly generated and maintained to reflect enhanced filtering capabilities (e.g., de-duplication and hate and profanity detection) and improved tooling.

C. Organization of Report

The remainder of this report is organized as follows. In Section II, we describe the data sources used in granite.13b's pre-training. In Section III, we describe the data processing steps we undertake with a focus on the governance steps we follow. In Section IV, we provide further details about the pre-training and fine-tuning algorithms, the computation involved, and the energy consumption we estimate. Section V presents the testing and evaluation framework along with quantitative comparisons to other models. In Section VI, we discuss our approach to understanding and mitigating socio-technical harms from the Granite models. Section VII provides a brief discussion of the usage policies and the socio-technical documentation of Granite models. Finally in Section VIII, we conclude with areas of future work and discussion.

II. DATA SOURCES

At the time of granite.13b's pre-training, IBM had curated 6.48 TB of data before pre-processing, 2.07 TB after pre-processing (detailed in Section III). All datasets were filtered English-text and code unstructured data files. There are no pre-defined labels or targets. All non-text artifacts (e.g., images, HTML tags, etc.) were removed.

Specifically, for the purposes of training granite.13b, 1 trillion tokens were generated from a total of 14 datasets. The individual datasets used in the training are described below.

- 1) *arXiv*: Over 1.8 million scientific paper pre-prints posted to arXiv.
- 2) *Common Crawl*: Open repository of web crawl data.
- 3) *DeepMind Mathematics*: Mathematical question and answer pairs data.
- 4) *Free Law*: Public-domain legal opinions from US federal and state courts.
- 5) *GitHub Clean*: Code data from CodeParrot covering a variety of coding languages.
- 6) *Hacker News*: News on computer science and entrepreneurship, taken between 2007-2018.
- 7) *OpenWeb Text*: Open-source version of OpenAI's Web Text corpus containing web pages through 2019.
- 8) *Project Gutenberg (PG-19)*: A repository of free e-books with focus on older works for which U.S. copyright has expired.

- 9) *Pubmed Central*: Biomedical and life sciences papers.
- 10) *SEC Filings*: 10-K/Q filings from the US Securities and Exchange Commission (SEC) for the years 1934-2022.
- 11) *Stack Exchange*: Anonymized set of all user-contributed content on the Stack Exchange network, a popular collection of websites centered around user-contributed questions and answers.
- 12) *USPTO*: US patents granted from 1975 to May 2023, excluding design patents.
- 13) *Webhose*: Unstructured web content converted into machine-readable data feeds acquired by IBM.
- 14) *Wikimedia*: Eight English Wikimedia projects (enwiki, enwikibooks, enwikinews, enwikiquette, enwikisource, enwikiversity, enwikivoyage, enwiktionary). containing extracted plain text from pages and articles.

III. DATA GOVERNANCE

As IBM is making Granite models available to customers to adapt to their own applications, we have invested heavily in a data governance process that evaluates datasets for governance, risk and compliance (GRC) criteria, including IBM's standard data clearance process, document quality checks, and other criteria. IBM has developed governance procedures for LLM pre-training datasets which are consistent with IBM AI Ethics principles and are guided by the IBM Corporate Legal Team. Best practices around LLM development is continually evolving with the ever-increasing understanding of AI models, their usage, and changing regulatory requirements, among other factors.

Addressing GRC criteria for data spans the lifecycle of training data, from data request to tokenization. An important objective for IBM is establishing an internal auditable link from a trained foundation model to the specific dataset version on which the model was trained, including information about each processing step performed prior to training. Summary statistics on IBM's curated pre-training dataset are provided in Fig. 2.

Data governance is organized into the following processes, corresponding to data lifecycle phases prior to model training:

- A. Data clearance and acquisition;
- B. Pre-processing; and
- C. Tokenization.

Each process is composed of sub-processes focusing on specific governance aspects. The remainder of this section describes each phase in detail.

A. Data Clearance and Acquisition

The data clearance process assures that no datasets are used to train IBM foundation models, including the Granite series, without careful consideration. Before data is added to IBM's curated pre-training dataset, it is submitted to the data clearance process and subject to technical, business,

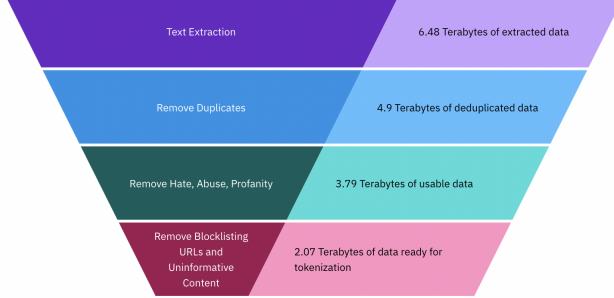


Fig. 2. Summary governance statistics on IBM’s curated pre-training dataset at the time of granite.13b’s training.

and governance review. The clearance request captures comprehensive information about a dataset such as a thorough description, the data owner, the intended use, geographic location, data classification, licensing information (if available), usage restrictions and sensitivity (e.g., personal information). Additional information includes who will have access to the data, and how the data will be acquired.

Once a dataset completes the review process, it is tagged for potential inclusion, its metadata is moved into a catalog of approved datasets, and it is downloaded and prepared for the subsequent pre-processing stages.

Remark. *IBM’s pre-training dataset currently addresses copyrighted material through selective use of URL blocklisting for websites known to disseminate pirated information. Examples of datasets that are block-listed include the Books3 dataset, which is specifically excluded from use due to concerns on the copyright status of the data and its use in model training.*

B. Pre-Processing Pipeline

Once data has been cleared and downloaded, it is prepared for model training through a variety of steps collectively referred to as the *pre-processing pipeline*. An overview of the pre-processing pipeline for this release of Granite models is depicted in Fig. 3 and is composed of the following steps:

- 1) Text extraction
- 2) De-duplication
- 3) Language identification
- 4) Sentence splitting
- 5) Hate, abuse and profanity annotation
- 6) Document quality annotation
- 7) URL block-listing annotation
- 8) Filtering
- 9) Tokenization.

Some pre-processing steps follow an annotation/filtering pattern, where documents or sentences are annotated first and filtered later during the filtering task according to threshold definitions. The completion of each pipeline step in the pipeline is logged. Logs are used to construct metadata reflecting the

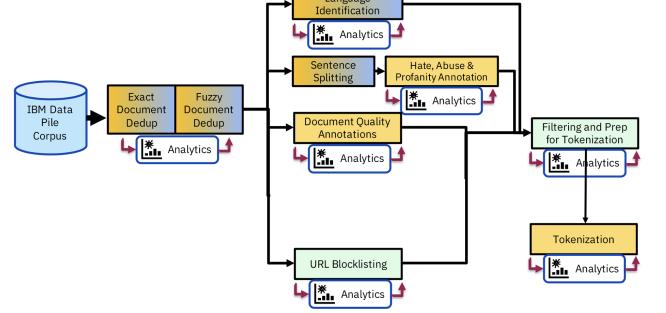


Fig. 3. IBM’s Data pre-processing pipeline.

exact pre-processing steps performed on a dataset, laying the basis for end-to-end traceability of the model lifecycle.

We now describe each step of the pre-processing pipeline in greater detail.

1) *Text Extraction:* Text extraction is the first step in the pipeline, and is used to extract language from various documents into a standardized format for further processing.

2) *Data De-Duplication:* Data de-duplication aims to identify and remove duplicate documents. De-duplication is performed on a per-dataset basis and is essential to ensuring the trained model does not learn artificial linguistic patterns due to repeated data in the dataset.

Two techniques are used: exact and fuzzy de-duplication, both of which use hash-based methods. As the name suggests, exact de-duplication removes exact duplicates among the documents in the dataset. Each document is hashed and documents with the same hash are fused to one. For example, if 50 documents in a dataset have the same hash, a single document will be used. Fuzzy de-duplication finds the Jaccard similarity between documents with locality sensitive hashing. If multiple updated snapshots of a dataset are downloaded, the exact de-duplication is performed across all snapshots.

3) *Language Identification:* Language identification is performed at a document level to detect the dominant language using the Watson Natural Language Processing (NLP) library [10]. The output of this task is an additional column in the parquet file containing a two letter ISO language code.

In the case of the Common Crawl dataset, language is already provided through folder names. The Watson NLP language identification algorithm is nevertheless run on Common Crawl documents, yielding two language classifications for these documents: Common Crawl and Watson NLP.

4) *Sentence Splitting:* Sentence splitting involves decomposing each document into its constituent sentences. Sentence splitting is key for hate, abuse, and profanity (HAP) annotation (to be discussed below) since HAP annotation is performed at a sentence level. As such, the sentence splitting stage must take place prior to the start of HAP annotation. Sentence splitting for the English language is performed using Watson NLP.

5) *Hate, Abuse and Profanity Annotation*: Data sources drawing from the open Internet, such as Common Crawl, inevitably contain abusive language. To reduce the possibility of Granite models producing profane content, each sentence in each document is assessed and scored as to its level of HAP content. The HAP detector is itself a language model trained by IBM and benchmarked against internal as well as public models such as OffensEval [11], AbusEval [12] and HatEval [13]. The IBM HAP detector performs comparably to HateBERT [14].

After a score is assigned to each sentence in the document, analytics are run over the sentences and scores to explore the distribution of annotations in each document with a HAP annotation. This serves both to determine the percentage of HAP sentences in a document as well as to determine threshold values used later during filtering.

6) *Document Quality*: Quality annotation aims to identify documents with low linguistic value using both heuristics and a classifier. The heuristics are derived from the Gopher Quality Filtering criteria [15]:

- total words: outside the range 50–100,000 words;
- average word length: outside the range 3–10 characters per word;
- symbol to word ratio: greater than 10%;
- bullet points ratio: greater than 90%;
- ellipsis line ratio: greater than 30%;
- alphabet words ratio: fewer than 80%;
- common English words: does not contain at least 2 from {the, be, to, of, and, that, have, with}.

The classifier assigns a perplexity score using the KenLM linear classifier pre-trained on Wikipedia documents [16], [17]. For any document, the model provides a score of the document’s similarity to a training corpus (i.e., Wikipedia).

These heuristics and classifiers output columns with quality scores that are added to the parquet file. These annotations form the basis for quality filtering during the filtering step.

7) *URL Block-Listing*: Block-listing identifies documents to be blocked from being added to IBM’s curated pre-training dataset. The block list is continuously maintained and includes URLs of known copyrighted material as well as block-listed sites such those contained in the 2022 Review of Notorious Markets for Counterfeiting and Piracy [18].

8) *Filtering*: Filtering occurs at the document level and is the last step before tokenization. It is here that annotations created in previous pre-processing steps are used to prevent documents from being used for tokenization. For example, documents are dropped which exceed HAP thresholds or do not meet a defined document quality. For the current English-only Granite models, the language identification annotations are used to filter out non-English documents.

C. Tokenization

Tokenization is the final pre-processing step prior to model training. For granite.13b, the cleaned and filtered text is

converted from a sequence of characters to a vector of tokens using the GPT-NeoX 20B tokenizer [6].

IV. TRAINING

In this section, we detail the training process for the decoder-only Granite models covering the algorithmic details of pre-training and fine-tuning, the computing involved, and an estimate of the carbon footprint.

A. Algorithmic Details

1) *Pre-Training*: We adopt most of the pre-training settings from [19]. Specifically, we use the standard decoder-only transformer architecture [20], Gaussian error linear unit (GELU) activation function [21], MultiQuery-Attention for inference efficiency [22], and learned absolute positional embeddings. We also adopt FlashAttention to speed up the training and reduce its memory footprint [23], allowing us to increase the context length to 8192 from the context length 2048 used by many existing LLMs.

The granite.13b base model is trained for 300K iterations, with a batch size of 4M tokens, for a total of 1 trillion tokens. We train using the Adam optimizer [24], with $\beta_1 = 0.9$, $\beta_2 = 0.95$, $\epsilon = 10^{-8}$, and a weight decay of 0.1. We use a cosine learning rate schedule, with warmup of 2000 steps, and decay final learning rate down from 3×10^{-4} to 3×10^{-5} . We pre-train models with a 3D-parallel layout using both tensor and pipeline parallelism including sequence parallelism to enable training with 8K context length.

2) *Supervised Fine-Tuning*: Pre-training teaches the LLM to continue generating text based on the input. However in practice, users often expect the LLM to treat the input as instructions to follow. To enable instruction following, we perform supervised fine-tuning (SFT) with a mixture of datasets from different sources. Each sample consists of a prompt and an answer. We use a cosine learning rate schedule with an initial learning rate of 2×10^{-5} , a weight decay of 0.1, a batch size of 128, and a sequence length of 8192 tokens. We perform SFT for 3 epochs to obtain the granite.13b.instruct model.

The SFT data includes a subset of the Flan Collection [25], 15K samples from Dolly [2], Anthropic’s human preference data about helpfulness and harmlessness [3], Instructv3 [26], and internal synthetic datasets specifically designed for summarization and dialogue tasks.

3) *Contrastive Fine-Tuning*: Contrastive fine-tuning (CFT) is an instruction fine-tuning approach based on unlikelihood-based training [27], which penalizes the probability of data points from a negative data distribution while simultaneously increasing the probability of data points from a positive data distribution (see Fig. 4). In other words, we discourage an LLM from generating misaligned responses (e.g. responses that are harmful) while encouraging aligned responses (e.g. responses that are helpful) for each training prompt.

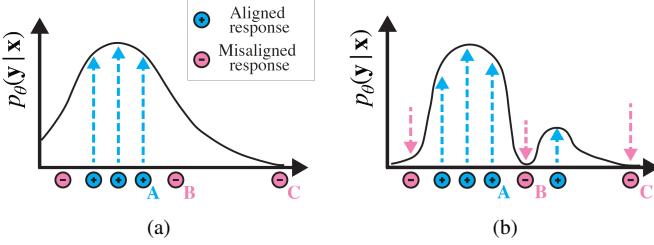


Fig. 4. Illustration for the resulting LLM distribution over responses y for a given prompt x . (a) SFT increases the likelihood of aligned responses but lacks control over misaligned responses. As a result, misaligned responses that closely resemble aligned ones can still have high likelihood. Examples A and B demonstrate this, where A is aligned and B is misaligned, yet they are similar. (b) Our approach, contrastive fine-tuning mitigates this by explicitly assigning low likelihood to misaligned responses.

CFT requires both responses to be paired with the same prompt in order for the model to determine which response is worse. However, many publicly available human demonstration datasets lack paired aligned and misaligned responses for the same prompt. As such, one may wonder: “*How can one obtain both aligned and misaligned responses?*” A straightforward approach to obtain this negative data distribution is to have humans write misaligned responses for each prompt. However, such an approach can be cost-prohibitive.

Thus in our work we propose to use a separate LLM to serve as a ‘negative persona’ by mimicking individuals who tend to respond in a misaligned (e.g. harmful or untruthful) manner. We achieve this by fine-tuning an LLM on widely available misaligned human demonstration datasets. Consequently, for a given dataset of prompts and aligned human demonstrations, responses on the prompts from this negative persona LLM form the misaligned responses and the paired aligned responses are the demonstrations themselves.

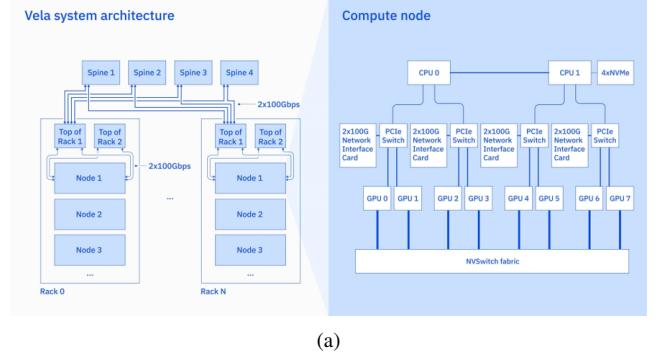
For granite.13b.chat, we use an early version of granite.13b.instruct as the separate LLM. The datasets for CFT are paired samples from Anthropic’s human preference data about helpfulness and harmlessness that have been filtered using the OpenAssist reward model [3], samples from Dolly [2], and samples from ProsocialDialog [4].

As a part of the CFT step, the granite.13b.chat was also trained to work with the following system prompt [3] in order to support Human-Agent based dialogue:

Below are a series of dialogues between various people and an AI assistant. The AI tries to be helpful, polite, honest, sophisticated, emotionally aware, and humble-but-knowledgeable. The assistant is happy to help with almost anything, and will do its best to understand exactly what is needed. It also tries to avoid giving false or misleading information, and it caveats when it isn’t entirely sure about the right answer. Moreover, the assistant prioritizes caution over usefulness, refusing to answer questions that it considers unsafe, immoral, unethical or dangerous.

Human:<prompt>

Assistant:



(a)

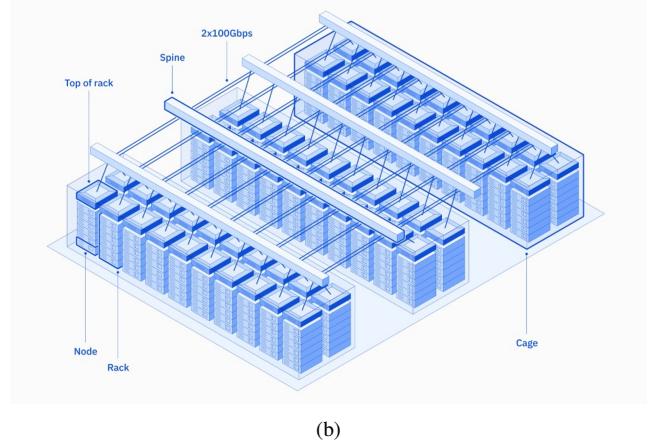


Fig. 5. An (a) architectural and (b) infrastructure diagram of the AI supercomputer Vela.

B. Compute

IBM’s primary computing infrastructure for training foundation models is the Vela AI supercomputer [28] (cf. diagram in Fig. 5). Vela uses a virtual machine-based approach for elasticity in resource allocation; with various optimizations, the ‘virtual machine tax’ is less than 5%. Each AI node has 8 Nvidia A100 GPU Cards, 96 vCPUs, 1.5 TB of DRAM and 4×3.2 TB NVMe drives. The nodes are interconnected via Ethernet. Each node has 2×100 Gbps Ethernet links. The Vela instance currently being used for model training is located in IBM Cloud’s Washington D.C. Data Center. Future Granite models are planned to be trained using Vela, however, the granite.13b base model was trained on older infrastructure before the Vela instance was fully stood up. Granite.13b used 256 A100 GPUs for 1056 hours and 120 TFLOPs.

C. Energy Consumption and Carbon Emissions

The methodology used to estimate the energy consumption and carbon emissions of the granite.13b base model is as follows. The carbon emissions $Carbon$ associated with a model M at a particular location L is given by:

$$Carbon(M, L) = E(M) \times PUE(L) \times CEF(L), \quad (1)$$

where $E(M)$ is the electricity consumption of the model M , $PUE(L)$ is the power usage effectiveness at the location L ,

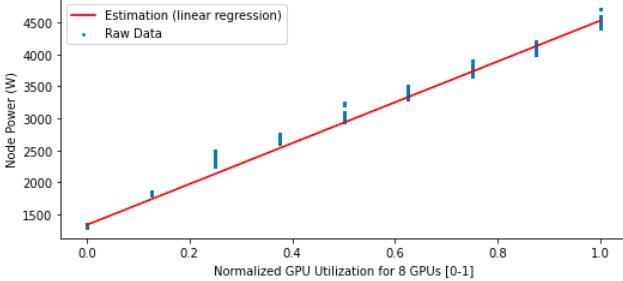


Fig. 6. Server (node) power vs. normalized GPU utilization.

and $CEF(L)$ is the carbon emission factor applicable for the location L .

The information technology (IT) electricity consumption $E(M)$ is estimated using the average GPU utilization rate for all the GPUs. It is a proxy to estimate the power that is used to train the AI model M since the GPU utilization is typically highly correlated with the node power, as shown in Fig. 6. Then, the estimated node power is multiplied by the training time and the number of GPUs used to calculate the total compute energy consumption E .

Power usage effectiveness $PUE(L)$ is given by the ratio of the total electricity consumed by the data center (aggregate consumption by the IT and support overhead infrastructure) to that consumed by the IT infrastructure. We calculate the location-based carbon emission factor $CEF(L)$ following the GHG Protocol’s Scope 2 Guidance [29].

Applying this estimation methodology to the granite.13b base model, we estimated 153074.3767 kWh energy consumption $E(M)$ and 0.12 kg/kWh carbon emission factor $CEF(L)$, yielding 22.2263995 tons of CO₂ equivalent $Carbon(M, L)$, which accounts for carbon dioxide and all other greenhouse gases, such as methane and nitrous oxide.

A number of mitigation strategies may be used to reduce the energy and carbon footprint. For example, the amount of resources used in training may be adjusted as a function of the availability of renewable energy, or the resources usage may be capped to not exceed certain energy usage or emissions limits.

V. TESTING AND EVALUATION

In this section, we describe the approach taken to test and evaluate the Granite models. We also provide empirical results along with comparisons to several other models that are of a similar capability level.

A. Foundation Model Evaluation Framework

We use a comprehensive foundation model evaluation framework (FM-eval) through the model’s development lifecycle. FM-eval is running on RedHat OpenShift¹ cluster with GPU

support, for efficient execution of evaluation benchmarks, in parallel and on multiple models. The automation framework can run any containerized evaluation framework or a wrapped external framework such as Eleuther AI’s Language Model Evaluation Harness (lm-eval) [30] or Stanford’s HELM (Holistic Evaluation Model) [31]. To allow easy addition of tasks, datasets and metrics to FM-eval, we developed Unitxt², an open-source Python library that provides a consistent interface and methodology for defining datasets, including the preprocessing required to convert raw datasets to the input required by LLMs, and the metrics used to evaluate the results.

Different types of tests are run during different phases of the lifecycle:

- 1) General knowledge benchmarks (during training)
- 2) HELM benchmarks (post-training)
- 3) Enterprise benchmarks (post-training)

These evaluations all leverage zero-shot and few-shot prompting. For clarity, zero-shot prompting uses a pre-existing LLM to generate text for a new task by only providing the instruction to execute the task in the prompt. In few-shot prompting, we provide multiple in-context examples, along with the task at hand, directly within the prompt. Both approaches allowed us to work with a single pre-trained model whose core parameters remained fixed.

The specific evaluations are detailed below.

1) *General Knowledge Benchmarks During Training*: The General Knowledge Benchmarks include a subset of existing benchmarks from lm-eval [30] and are used as light-weight tests run after every 100 billion tokens during training to validate model knowledge is advancing as training progresses.

Specifically, the following 12 datasets (organized by task) from lm-eval are:

- question answering for several domains (boolq, openbookqa, piqa, sciq);
- sentence completion (lambada)
- commonsense reasoning (arc_easy, arc_challenge, copa, hellaswag, winogrande);
- reading comprehension (race)
- multidisciplinary multiple-choice collection (mmlu);

In our evaluation framework these benchmarks are run in both the zero-shot and few-shot setting.

2) *HELM*: After pre-training is complete, a more comprehensive assessments relies on Stanford’s Holistic Evaluation of Language Models (HELM) Benchmark [31]. To evaluate our model, we use the 16 “core scenarios”, consisting of a variety of tasks including question answering, information retrieval, summarization, sentiment analysis, and text classification [32]–[43], on which all HELM LLMs are evaluated and compared using their results and mean win rate (MWR).

3) *Enterprise Evaluation Benchmarks*: After training completes, we further evaluate our models on IBM-curated enterprise benchmarks to test our models performance in domains

¹<https://www.redhat.com/en/technologies/cloud-computing/openshift>

²<https://github.com/IBM/unitxt>

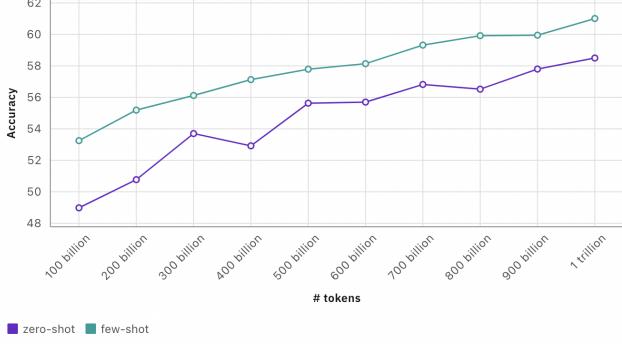


Fig. 7. Granite.13b General Knowledge Performance during Training.

highly relevant to our customers. With this in mind, IBM curated 11 publicly available finance benchmarks for evaluating models in the financial domain, summarized in Table I. The data source-provided train and test splits are used in the evaluation whenever possible. Model performance is reported based on test examples. If the test labels are not publicly available, model performance is reported on the validation set. If the train and test splits do not exist in the data source, 20% of the data is selected as test split and the rest is used as the train split.

All few-shot context examples are sampled from the training set. The number of few-shot examples provided to the model depends on the task, which is provided in Table I. For the current evaluation, all the models used the same parameters and the same standard prompt (see the techniques of few-shot-prompting and zero-shot-prompting and examples of prompts³), without task description, chain-of-thought prompting [44], or system prompts in place. For Financial Phrasesbank, News Headline and FiQA SA, the prompts were taken from BloombergGPT [45].

B. Granite Model Evaluation and Comparison

It should be noted that while future versions of granite.13b will be trained on upwards of 2T tokens, this initial release of granite.13b only saw 1T tokens during training, making all of these evaluations preliminary.

1) General Knowledge Benchmarks During Training: In this section, we leverage the lighter-weight General Knowledge Benchmarks to assess a series of snapshots of the granite.13b base model taken every 100B tokens during training along with the fine-tuned granite.13b.instruct and granite.13b.chat variants. As visualized in Fig. 7 and further detailed in Table II, progressively training on each 100B tokens steadily improved General Knowledge as expected with further boosts in performance achieved in both fine-tuned variants of granite.13b. Note system prompts were not used for this evaluation of granite.13b.chat.

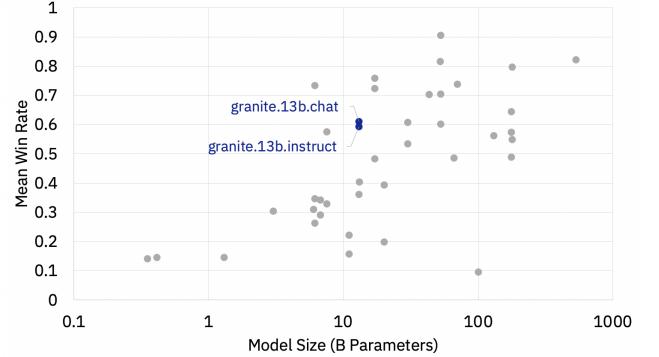


Fig. 8. Model Performance on HELM Tasks vs Model Size.

2) HELM Benchmarks: In this study, we comprehensively assess our models across HELM’s 16 core scenarios, comparing them to all other models in the v0.2.3 release⁵.

Our evaluation follows the default two-tiered process as suggested by HELM: first, we evaluate each model on individual evaluation datasets, then we aggregate these results into scenarios (as detailed in Table III). To facilitate a fair comparison of LLMs, we employ HELM’s Mean Win Rate (MWR) metric across scenarios for model ranking.

Figure 8 illustrates the positioning of granite.13b.instruct, granite.13b.chat, and all v0.2.3 models along the model size MWR axes. Note, models where exact model size is not published by the LLM-provider are excluded from this visualization.

The figure highlights that the granite models strike a desirable balance between model size and HELM performance. Granite.13b.chat and granite.13b.instruct are respectively ranked 15 and 18, out of all models evaluated. Further, the granite.13b.chat and granite.13b.instruct models were respectively the top-2 and top-3 models evaluated under 17B parameters in size. Only Cohere Command beta (6.1B) exceeded their performance in this size category. These results also hold for the other aspects evaluated by HELM, such as robustness, and fairness. In calibration granite.13b.chat and granite.13b.instruct are ranked 9 and 28, respectively.

3) Enterprise Benchmarks: This evaluation is conducted by augmenting HELM’s framework to encompass 11 publicly available task datasets from the financial services domain. Baseline models are selected based on model size, type of training data, accessibility, and model tuning. To be specific, granite models are compared with GPT-NeoX-20B [6], Pythia-12B-sft-4 [56] are FLAN-UL2 [57] that are among the best performing open-sourced models under 50 billion parameters. In addition, Pythia-12B-sft-4 [56] is an SFT model. The baseline models also include the state-of-the-art available models, LLaMA2 [58], with 7 billion to 13 billion parameters.

Table IV presents the detailed performance scores of the models on the 11 financial tasks. An asterisk is given next to the Llama2 models, as these models have seen 2T tokens

³<https://www.promptingguide.ai/techniques/fewshot>

⁵https://crfm.stanford.edu/helm/latest/?group=core_scenarios

TABLE I
FINANCE BENCHMARKS OVERVIEW

| Task | Task Description | Dataset | Dataset Description | N-shot Prompt | Metric |
|--------------------------|----------------------------|---|--|---------------|----------------------|
| Sentiment Classification | 3 classes | Financial Phrasebank [46] | Financial news categorised by sentiment | 5-shot | Weighted F1 |
| | 2 classes | Earnings Call Transcripts [47] | Earnings call transcripts, the related stock prices and the sector index in terms of volume | 5-shot | Weighted F1 |
| Classification | 9 classes | News Headline [48] | The gold commodity news annotated into various dimensions | 5-shot | Weighted F1 |
| Named Entity Recognition | 4 numerical entities | Credit Risk Assessment (NER) [49] | Eight financial agreements (totalling 54,256 words) from SEC filings were manually annotated for entity types: location, organization person and miscellaneous | 20-shot | Entity F-1 |
| | 4522 numerical entities | KPI-Edgar [50] | A dataset for Joint Named Entity Recognition and Relation Extraction building on financial reports uploaded to the Electronic Data Gathering, Analysis, and Retrieval (EDGAR) system, where the main objective is to extract Key Performance Indicators (KPIs) from financial documents and link them to their numerical values and other attributes | 20-shot | Modified Adjusted F1 |
| | 139 numerical entities | FiNER-139 [51] | 1.1M sentences annotated with extensive Business Reporting Language (XBRL) tags extracted from annual and quarterly reports of publicly-traded companies in the US, focusing on numeric tokens, with the correct tag depending mostly on context, not the token itself. | 10-shot | Entity F1 |
| Question Answering | Document relevance ranking | Opinion-based QA (FiQA) [52] | Text documents from different financial data sources (microblogs, reports, news) for ranking document relevance based on opinionated questions, targeting mined opinions and their respective entities, aspects, sentiment polarity and opinion holder. | 5-shot | RR@10 |
| | 3 classes | Sentiment Analysis (FiQA SA) [52] | Text instances in the financial domain (microblog message, news statement or headline) for detecting the target aspects which are mentioned in the text (from a pre-defined list of aspect classes) and predict the sentiment score for each of the mentioned targets. | 5-shot | Weighted F1 |
| | Ranking | Insurance QA [53] | Questions from real world users and answers with high quality composed by professionals with deep domain knowledge collected from the website Insurance Library ⁴ | 5-shot | RR@10 |
| | Exact value match | Chain of Numeric Reasoning (ConvFinQA) [54] | Multi-turn conversational finance question answering data for exploring the chain of numerical reasoning | 1-shot | Accuracy |
| Summarization | Long documents | Financial text summarization (EDT) [55] | 303893 news articles range from March 2020 to May 2021 for abstractive text summarization | 5-shot | Rouge-L |

TABLE II
GRANITE.13B GENERAL KNOWLEDGE PERFORMANCE DURING TRAINING

| Model | Tokens (B) | Avg Accuracy (Zero-Shot) | Avg Accuracy (Few-Shot) |
|-----------------------------|-------------|--------------------------|-------------------------|
| granite.13b (base) | 100 | 49.0 | 53.3 |
| granite.13b (base) | 200 | 50.8 | 55.2 |
| granite.13b (base) | 300 | 53.7 | 56.1 |
| granite.13b (base) | 400 | 52.9 | 57.1 |
| granite.13b (base) | 500 | 55.6 | 57.8 |
| granite.13b (base) | 600 | 55.7 | 58.1 |
| granite.13b (base) | 700 | 56.8 | 59.3 |
| granite.13b (base) | 800 | 56.5 | 59.9 |
| granite.13b (base) | 900 | 57.8 | 60.0 |
| granite.13b (base) | 1000 | 58.5 | 61.0 |
| granite.13b.instruct | 1000 | 59.3 | 61.5 |
| granite.13b.chat | 1000 | 61.2 | 62.6 |

of pre-training data, imparting significant advantage to the models. All other models evaluated, including granite, have seen 1T tokens of training data except for the pythia model which has only seen 300B tokens during pre-training. Despite having been trained on half the amount of data as the Llama2 models, the Granite models are competitive across each of the tasks, often outperforming Llama2. This bodes well for future

planned versions of granite models, which will be trained on 2T+ tokens of pre-training data.

VI. SOCIO-TECHNICAL HARMS AND RISKS

Numerous potential socio-technical harms and risks of LLMs have been identified in recent years, including misinformation, hallucination, lack of faithfulness or factuality, leakage of private information, plagiarism or inclusion of copyrighted content, hate speech, toxicity, human-computer interaction harms such as bullying and gaslighting, malicious uses, and adversarial attacks [59], [60].

In Table V, we present the catalogue of risks compiled by the IBM AI Ethics Board, a central, cross-disciplinary body that defines the AI ethics vision and strategy with the objective of supporting a culture of ethical, responsible, and trustworthy AI throughout the IBM Corporation [61], [62]. The table is organized across several dimensions [63]:

- Whether the risk is from the data or other inputs to the foundation model, from the generated output of the foundation model, or from other concerns.

TABLE III
HELM RESULTS PER SUB SCENARIO IN THE CORE-SCENARIOS.

| Model (Metric) | MMLU (EM) | BoolQ (EM) | NarrativeQA (F1) | NaturalQuestions closed-book (F1) | NaturalQuestions open-book (F1) | QuAC (F1) | HellaSwag (EM) | OpenbookQA (EM) | TruthfulQA (EM) | MS MARCO (RR@10) | MS MARCO (TREC) (NDCG@10) | CNN/DailyMail (ROUGE-2) | XSUM (ROUGE-2) | IMDB (EM) | CivilComments (EM) | RAFT (EM) |
|----------------------|-----------|------------|------------------|-----------------------------------|---------------------------------|-----------|----------------|-----------------|-----------------|------------------|---------------------------|-------------------------|----------------|-----------|--------------------|-----------|
| granite.13b.instruct | 0.377 | 0.809 | 0.668 | 0.188 | 0.659 | 0.373 | 0.338 | 0.296 | 0.203 | 0.431 | 0.638 | 0.135 | 0.11 | 0.953 | 0.637 | 0.693 |
| granite.13b.chat | 0.378 | 0.776 | 0.698 | 0.212 | 0.684 | 0.391 | 0.305 | 0.276 | 0.208 | 0.396 | 0.634 | 0.14 | 0.115 | 0.948 | 0.6 | 0.709 |

TABLE IV
FINANCE BENCHMARK EVALUATION RESULTS PER TASK.

| | Financial Phrase-bank | Earnings Call Transcripts | News Headline | Credit Risk Assessment | KPI-Edgar | FiNER-139 | FiQA - Opinion | Insurance QA | FiQA SA | ConFinQA | Summarization |
|----------------------|-----------------------|---------------------------|---------------|------------------------|--------------|--------------|----------------|--------------|--------------|--------------|---------------|
| Metrics | Weighted F1 | Weighted F1 | Weighted F1 | Entity F1 | Adj F1 | Entity F1 | RR @10 | RR@10 | Weighted F1 | Accuracy | Rough-L |
| granite.13b (base) | 0.306 | 0.443 | 0.811 | 0.477 | 0.344 | 0.699 | 0.398 | 0.365 | 0.780 | 0.365 | 0.173 |
| granite.13b.instruct | 0.590 | 0.443 | 0.764 | 0.407 | 0.281 | 0.699 | 0.489 | 0.979 | 0.590 | 0.346 | 0.323 |
| granite.13b.chat | 0.714 | 0.443 | 0.779 | 0.361 | 0.290 | 0.746 | 0.486 | 0.990 | 0.758 | 0.334 | 0.376 |
| llama2.7b* | 0.244 | 0.486 | 0.752 | 0.408 | 0.419 | 0.660 | 0.548 | 0.365 | 0.744 | 0.233 | 0.195 |
| llama2.7b.chat* | 0.758 | 0.677 | 0.829 | 0.458 | 0.450 | 0.626 | 0.554 | 0.346 | 0.693 | 0.254 | 0.345 |
| llama2.13b* | 0.378 | 0.410 | 0.584 | 0.467 | 0.463 | 0.689 | 0.409 | 0.365 | 0.800 | 0.226 | 0.252 |
| llama2.13b.chat* | 0.608 | 0.572 | 0.744 | 0.445 | 0.538 | 0.671 | 0.532 | 0.346 | 0.849 | 0.261 | 0.269 |
| gpt-neox-20b | 0.561 | 0.318 | 0.630 | 0.469 | 0.308 | 0.774 | 0.446 | 0.865 | 0.771 | 0.266 | 0.205 |
| flan-ul2 | 0.240 | 0.318 | 0.829 | 0.394 | 0.011 | 0.446 | 0.695 | 0.708 | 0.811 | 0.254 | 0.310 |
| oasst-sft-pythia-12b | 0.536 | 0.318 | 0.579 | 0.105 | 0.177 | 0.514 | 0.480 | 0.802 | 0.752 | 0.217 | 0.086 |

- Whether the risk arises in the training/tuning of the model, during inference, or in broader considerations such as governance, legal compliance, or societal impact.
- What higher-level grouping the risk falls under, e.g. fairness, robustness, intellectual property, and misuse.
- Whether the risk is new or amplified. ‘Traditional’ risks are present in earlier forms of AI models and continue to be present in foundation models. ‘Amplified’ risks are known from earlier forms of AI models but are intensified by foundation models due to their generative capabilities. ‘New’ risks are emerging risks, intrinsic to foundation models due to their generative capabilities.

As part of creating and releasing the granite.13b.instruct and granite.13b.chat models, we have addressed some of the risks as follows. The data governance processes of the IBM’s pre-training dataset, including the block-listing and filtering of hate, abuse and profanity have mitigated some intellectual property and misuse risks. Toward fairness, an additional component of the data pre-processing pipeline not described in Section III is annotating documents by religion, gender, race, stigma, age, and political ideology. We have created keyword lists for these dimensions and use keyword matching to annotate sentences. The annotations may be used to identify under-represented and over-represented groups. We have not been overly aggressive in HAP filtering and have not filtered with respect to groups because it would prevent us from having training data that reclaims slurs and positively describes marginalized identities, and might skew the pre-training dataset in other unintended ways [64].

Through fine-tuning, we have encouraged prosocial and less harmful model behavior with the aim to mitigate certain aspects of misuse and value alignment risks. However, one of the biggest socio-technical risks would be our own hubris

to believe that the datasets we used for fine-tuning (or other existing datasets we could use in the future) are aligned to the needs, wants, and desires of the peoples and organizations that will be deploying the Granite models to meet their own goals. Every enterprise has its own regulations to conform to, whether they come from laws, social norms, industry standards, market demands, or architectural requirements [65]; we believe that enterprises should be empowered to personalize their models according to their own values (within bounds) [66], e.g. using tools in the watsonx platform.

In addition, through FM-eval, we have tested the Granite models on benchmark datasets that cover several risk dimensions. However, evaluating on benchmarks is a limited approach for revealing socio-technical harms [67]. After enterprises have further aligned the Granite models to their own values, they should enlist a red team with members of varying socio-cultural and lived experience to find additional harms and undesirable LLM behaviors within the context of a precise use case [68].

VII. USAGE POLICIES AND DOCUMENTATION

A. Machine-Generated Content

IBM’s licensing terms and conditions govern downstream applications and services that use IBM models.

In addition, IBM is setting up an Acceptable Use Provision (AUP) that states guidelines and practices that the users of IBM models are required to follow as they develop and deliver downstream applications and services.

The AUP provides acceptable use of AI Models and confers to IBM the right to terminate the license to these models if necessary.

TABLE V
SOCIO-TECHNICAL HARMS AND RISKS

| Source | Phase | Group | Risk | Indicator |
|--------|---------------------|-------------------------|--|-------------|
| Input | Training and Tuning | Fairness | Bias | Amplified |
| Input | Training and Tuning | Robustness | False samples | Traditional |
| Input | Training and Tuning | Value Alignment | Undesirable output for retraining purposes | New |
| Input | Training and Tuning | Data Laws | Legal restrictions on moving or using data | Traditional |
| Input | Training and Tuning | Intellectual Property | Copyright and other IP issues with content | Amplified |
| Input | Training and Tuning | Transparency | Disclose data collected, who has access, how stored, how it will be used | Amplified |
| Input | Training and Tuning | Privacy | Inclusion or presence of SPI or PII | Traditional |
| Input | Training and Tuning | Privacy | Provide data subject rights (e.g., opt-out) | Amplified |
| Input | Inference | Privacy | Disclose PII or SPI as part of prompt to model | New |
| Input | Inference | Intellectual Property | Disclose copyright or other IP information as part of prompt to model | New |
| Input | Inference | Robustness | Vulnerabilities to adversarial attacks like evasion (create incorrect model output by modifying data sent to train model) | Amplified |
| Input | Inference | Robustness | Vulnerabilities to adversarial attacks like prompt injection (force different output), prompt leaking (disclose system prompt), or jailbreaking (avoid guardrails) | New |
| Output | Inference | Fairness | Bias in generated content | New |
| Output | Inference | Fairness | Performance disparity across individuals or groups | Traditional |
| Output | Inference | Intellectual property | Copyright infringement, compliance with open source license agreements | New |
| Output | Inference | Value alignment | Hallucination (generation of false content) | New |
| Output | Inference | Value alignment | Toxic, hateful, abusive, and aggressive output | New |
| Output | Inference | Misuse | Spread disinformation (deliberate creation of misleading information) | Amplified |
| Output | Inference | Misuse | Generate toxic, hateful, abusive, and aggressive content | New |
| Output | Inference | Misuse | Nonconsensual use of people's likeness (deepfakes) | Amplified |
| Output | Inference | Misuse | Dangerous use (e.g., creating plans to develop weapons or malware) | New |
| Output | Inference | Misuse | Deceptive use of generated content (e.g., intentional nondisclosure of AI generated content) | New |
| Output | Inference | Harmful code generation | Execution of harmful generated code | New |
| Output | Inference | Privacy | Expose PI or SPI in generated content | New |
| Output | Inference | Explainability | Challenges in explaining the generated output | New |
| Output | Inference | Traceability | Challenges in identifying source and facts for generated output | New |
| Other | Governance | Transparency | Document data and model details, purpose, potential use and harms | Traditional |
| Other | Governance | Accountability | Identify responsibility for misaligned output along AI lifecycle and value chain | Amplified |
| Other | Legal compliance | Intellectual property | Determine creator of downstream models | New |
| Other | Legal compliance | Intellectual property | Determine creator of open source foundation models | New |
| Other | Legal compliance | Intellectual property | Determine owner of AI-generated content | New |
| Other | Legal compliance | Intellectual property | Uncertainty about IP rights related to generated content | New |
| Other | Legal compliance | Legal uncertainty | Determine downstream obligations | Amplified |
| Other | Societal impact | Impact on jobs | Human displacement (AI induced job loss) | Amplified |
| Other | Societal Impact | Human dignity | Human exploitation (ghost work in training), poor working conditions, lack of healthcare, unfair compensation | Amplified |
| Other | Societal Impact | Environment | Increased carbon emission (high energy requirements for training and operation) | Amplified |
| Other | Societal Impact | Diversity and inclusion | Homogenizing culture and thoughts | New |
| Other | Societal Impact | Human agency | Misinformation and disinformation generated by foundation models | Amplified |
| Other | Societal Impact | Impact on education | Bypass learning process, plagiarism | New |

B. European Union-Specific Controls

The licensing terms and conditions to IBM Models are augmented with an Acceptable Use Policy (AUP) that states guidelines and practices that are specific to deployments of downstream applications and services in specific Countries.

C. Downstream Documentation

For downstream usage of its pre-trained models, IBM makes available the following documentation:

- Terms and Conditions
- Product documentation
- Technical reports, such as this report

Together, this information is designed so that not only IBM complies with legal and ethical requirements, but also the users of these models can comply with their own obligations.

1) Terms and Conditions: The latest Terms and Conditions for the watsonx platform can be found at <https://www.ibm.com/support/customer/csolt/terms/?id=i126-6883>.

2) Product documentation: The IBM Granite models are currently available through IBM's watsonx platform. As part of watsonx, each Granite model is accompanied by a model card that details key facts and provenance of the model.

VIII. CONCLUSION

In this technical report, we have presented IBM's Granite family of foundation models designed for enterprise generative AI applications. IBM's ethical and governance frameworks provide the context within which these models are created and made available. Aligned with IBM's commitment to transparent and responsible AI, we have presented descriptions of exact datasets, pre-processing steps, training infrastructure, energy consumption, and testing/evaluation methodologies used throughout the model development lifecycle.

We are continuing to develop the Granite series in several directions. Whereas this initial Granite release only supports English, future models will be trained on multiple natural languages. Alongside, HAP annotation is being refined and expanded for additional languages. Furthermore, Granite models for other modalities such as code as well as industry-specific content are being developed. On the model safety evaluation front, we are developing a comprehensive red-teaming framework. The adversarial prompts will test the models across a variety of domains, including (but not limited to) HAP, bias and stigma, factual correctness and harmful topics.

We are continuing to develop additional data annotations for IBM's curated pre-training dataset, such as scoring documents for their inclusion of personally-identifiable information and for their conversationality [69], [70]. We are working toward instrumenting our compute infrastructure to obtain precise rather than estimated measurement of energy and carbon

footprints [71]. Finally, we are exploring the application of various methods for mitigating unwanted biases [72]–[74].

REFERENCES

- [1] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le, "Finetuned language models are zero-shot learners," 2021.
- [2] M. Conover, M. Hayes, A. Mathur, J. Xie, J. Wan, S. Shah, A. Ghodsi, P. Wendell, M. Zaharia, and R. Xin, "Free Dolly: Introducing the world's first truly open instruction-tuned LLM," <https://www.databricks.com/blog/2023/04/12/dolly-first-open-commercially-viable-instruction-tuned-lm>, Apr. 2023.
- [3] Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli, T. Henighan *et al.*, "Training a helpful and harmless assistant with reinforcement learning from human feedback," *arXiv preprint arXiv:2204.05862*, 2022.
- [4] H. Kim, Y. Yu, L. Jiang, X. Lu, D. Khashabi, G. Kim, Y. Choi, and M. Sap, "ProsocialDialog: A prosocial backbone for conversational agents," in *Proc. Conf. Empir. Meth. Nat. Lang. Proc.*, Abu Dhabi, United Arab Emirates, Dec. 2022, pp. 4005–4029.
- [5] D. D. Cox, "Introducing the technology behind watsonx.ai, IBM's AI and data platform for enterprise," <https://www.ibm.com/blog/introducing-the-technology-behind-watsonx-ai>, May 2023.
- [6] S. Black, S. Biderman, E. Hallahan, Q. Anthony, L. Gao, L. Golding, H. He, C. Leahy, K. McDonell, J. Phang, M. Pieler, U. S. Prashanth, S. Purohit, L. Reynolds, J. Tow, B. Wang, and S. Weinbach, "Gpt-neox-20b: An open-source autoregressive language model," 2022.
- [7] L. Gao, S. Biderman, S. Black, L. Golding, T. Hoppe, C. Foster, J. Phang, H. He, A. Thite, N. Nabeshima, S. Presser, and C. Leahy, "The Pile: An 800gb dataset of diverse text for language modeling," *arXiv preprint arXiv:2101.00027*, 2020.
- [8] <https://huggingface.co/datasets/c4>.
- [9] K. Schaul, S. Y. Chen, and N. Tiku, "Inside the secret list of websites that make AI like ChatGPT sound smart," *Washington Post*, Apr. 2023.
- [10] IBM Corporation. Watson Natural Language Processing library. [Online]. Available: <https://dataplatform.cloud.ibm.com/docs/content/wsj/analyze-data/watson-nlp.html?context=cpdaas>
- [11] M. Zampieri, S. Malmasi, P. Nakov, S. Rosenthal, N. Farra, and R. Kumar, "SemEval-2019 task 6: Identifying and categorizing offensive language in social media (offenseval)," *CoRR*, vol. abs/1903.08983, 2019. [Online]. Available: <http://arxiv.org/abs/1903.08983>
- [12] T. Caselli, V. Basile, J. Mitrović, I. Kartoziya, and M. Granitzer, "I feel offended, don't be abusive! implicit/explicit messages in offensive and abusive language," in *Proceedings of the Twelfth Language Resources and Evaluation Conference*. Marseille, France: European Language Resources Association, May 2020, pp. 6193–6202. [Online]. Available: <https://aclanthology.org/2020.lrec-1.760>
- [13] A. Capozzi, M. Lai, V. Basile, C. Musto, M. Polignano, F. Poletto, M. Sanguinetti, C. Bosco, V. Patti, G. Ruffo, G. Semeraro, and M. Stranisci, "Computational linguistics against hate: Hate speech detection and visualization on social media in the "contro l'odio" project," 11 2019.
- [14] T. Caselli, V. Basile, J. Mitrovic, and M. Granitzer, "Hatebert: Retraining BERT for abusive language detection in english," *CoRR*, vol. abs/2010.12472, 2020. [Online]. Available: <https://arxiv.org/abs/2010.12472>
- [15] J. W. Rae *et al.*, "Scaling Language Models: Methods, Analysis & Insights from Training Gopher," 2022. [Online]. Available: <https://arxiv.org/abs/2112.11446>
- [16] Kenneth Heafield. (2011) KenLM: Faster and smaller language model queries. [Online]. Available: <https://kheafield.com/papers/avenue/kenlm.pdf>
- [17] kenlm GitHub source code repository. [Online]. Available: <https://github.com/kpu/kenlm>
- [18] Office of the United States Trade Representative (USTR). (2022) 2022 Review of Notorious Markets for Counterfeiting and Piracy. [Online]. Available: [https://ustr.gov/sites/default/files/2023-01/2022%20Notorious%20Markets%20List%20\(final\).pdf](https://ustr.gov/sites/default/files/2023-01/2022%20Notorious%20Markets%20List%20(final).pdf)
- [19] R. Li, L. B. Allal, Y. Zi, N. Muennighoff, D. Kocetkov, C. Mou, M. Marone, C. Akiki, J. Li, J. Chiu *et al.*, "Starcoder: may the source be with you?" *arXiv preprint arXiv:2305.06161*, 2023.
- [20] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

- [21] D. Hendrycks and K. Gimpel, "Gaussian error linear units (gelus)," *arXiv preprint arXiv:1606.08415*, 2016.
- [22] N. Shazeer, "Fast transformer decoding: One write-head is all you need," *arXiv preprint arXiv:1911.02150*, 2019.
- [23] T. Dao, D. Fu, S. Ermon, A. Rudra, and C. Ré, "Flashattention: Fast and memory-efficient exact attention with io-awareness," *Advances in Neural Information Processing Systems*, vol. 35, pp. 16 344–16 359, 2022.
- [24] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [25] S. Longpre, L. Hou, T. Vu, A. Webson, H. W. Chung, Y. Tay, D. Zhou, Q. V. Le, B. Zoph, J. Wei *et al.*, "The flan collection: Designing data and methods for effective instruction tuning," *arXiv preprint arXiv:2301.13688*, 2023.
- [26] (2023) Mpt-30b: Raising the bar for open-source foundation models. [Online]. Available: <https://www.mosaicml.com/blog/mpt-30b>
- [27] S. Welleck, I. Kulikov, S. Roller, E. Dinan, K. Cho, and J. Weston, "Neural text generation with unlikelihood training," *arXiv preprint arXiv:1908.04319*, 2019.
- [28] T. Gershon, S. Seelam, J. Jubran, E. Gampel, and D. Thorstensen, "Why we built an AI supercomputer in the cloud," <https://research.ibm.com/blog/AI-supercomputer-Vela-GPU-cluster>, Feb. 2023.
- [29] Mary Sotos. (2015) GHG Protocol Scope 2 Guidance. [Online]. Available: https://ghgprotocol.org/sites/default/files/ghgp/standards/Scope%202%20Guidance_Final_0.pdf
- [30] L. Gao, J. Tow, S. Biderman, S. Black, A. DiPofi, C. Foster, L. Golding, J. Hsu, K. McDonell, N. Muennighoff, J. Phang, L. Reynolds, E. Tang, A. Thite, B. Wang, K. Wang, and A. Zou, "A framework for few-shot language model evaluation," Sep. 2021. [Online]. Available: <https://doi.org/10.5281/zendodo.5371628>
- [31] P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang, D. Narayanan, Y. Wu, A. Kumar, B. Newman, B. Yuan, B. Yan, C. Zhang, C. Cosgrove, C. D. Manning, C. Ré, D. Acosta-Navas, D. A. Hudson, E. Zelikman, E. Durmus, F. Ladzhak, F. Rong, H. Ren, H. Yao, J. Wang, K. Santhanam, L. Orr, L. Zheng, M. Yuksekgonul, M. Suzgun, N. Kim, N. Guha, N. Chatterji, O. Khattab, P. Henderson, Q. Huang, R. Chi, S. M. Xie, S. Santurkar, S. Ganguli, T. Hashimoto, T. Icard, T. Zhang, V. Chaudhary, W. Wang, X. Li, Y. Mai, Y. Zhang, and Y. Koreeda, "Holistic evaluation of language models," 2022.
- [32] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. X. Song, and J. Steinhardt, "Measuring massive multitask language understanding," *ArXiv*, vol. abs/2009.03300, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:221516475>
- [33] C. Clark, K. Lee, M.-W. Chang, T. Kwiatkowski, M. Collins, and K. Toutanova, "Boolq: Exploring the surprising difficulty of natural yes/no questions," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019, pp. 2924–2936.
- [34] T. Kočiský, J. Schwarz, P. Blunsom, C. Dyer, K. M. Hermann, G. Melis, and E. Grefenstette, "The NarrativeQA reading comprehension challenge," *Transactions of the Association for Computational Linguistics*, vol. 6, pp. 317–328, 2018. [Online]. Available: <https://aclanthology.org/Q18-1023>
- [35] R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi, "HellaSwag: Can a machine really finish your sentence?" in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Florence, Italy: Association for Computational Linguistics, Jul. 2019, pp. 4791–4800. [Online]. Available: <https://aclanthology.org/P19-1472>
- [36] T. Mihaylov, P. Clark, T. Khot, and A. Sabharwal, "Can a suit of armor conduct electricity? a new dataset for open book question answering," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Brussels, Belgium: Association for Computational Linguistics, Oct.-Nov. 2018, pp. 2381–2391. [Online]. Available: <https://aclanthology.org/D18-1260>
- [37] T. Kwiatkowski, J. Palomaki, O. Redfield, M. Collins, A. Parikh, C. Alberti, D. Epstein, I. Polosukhin, J. Devlin, K. Lee, K. Toutanova, L. Jones, M. Kelcey, M.-W. Chang, A. M. Dai, J. Uszkoreit, Q. Le, and S. Petrov, "Natural questions: A benchmark for question answering research," *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 452–466, 2019. [Online]. Available: <https://aclanthology.org/Q19-1026>
- [38] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts, "Learning word vectors for sentiment analysis," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*. Portland, Oregon, USA: Association for Computational Linguistics, June 2011, pp. 142–150. [Online]. Available: <http://www.aclweb.org/anthology/P11-1015>
- [39] D. F. Campos, T. Nguyen, M. Rosenberg, X. Song, J. Gao, S. Tiwary, R. Majumder, L. Deng, and B. Mitra, "Ms marco: A human generated machine reading comprehension dataset," *ArXiv*, vol. abs/1611.09268, 2016. [Online]. Available: <https://api.semanticscholar.org/CorpusID:1289517>
- [40] S. Narayan, S. B. Cohen, and M. Lapata, "Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Brussels, Belgium: Association for Computational Linguistics, Oct.-Nov. 2018, pp. 1797–1807. [Online]. Available: <https://aclanthology.org/D18-1206>
- [41] S. Lin, J. Hilton, and O. Evans, "TruthfulQA: Measuring how models mimic human falsehoods," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 3214–3252. [Online]. Available: <https://aclanthology.org/2022.acl-long.229>
- [42] A. See, P. J. Liu, and C. D. Manning, "Get to the point: Summarization with pointer-generator networks," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vancouver, Canada: Association for Computational Linguistics, Jul. 2017, pp. 1073–1083. [Online]. Available: <https://aclanthology.org/P17-1099>
- [43] E. Choi, H. He, M. Iyyer, M. Yatskar, W.-t. Yih, Y. Choi, P. Liang, and L. Zettlemoyer, "QuAC: Question answering in context," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Brussels, Belgium: Association for Computational Linguistics, Oct.-Nov. 2018, pp. 2174–2184. [Online]. Available: <https://aclanthology.org/D18-1241>
- [44] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," 2023.
- [45] S. Wu, O. Irsøy, S. Lu, V. Dabrowski, M. Dredze, S. Gehrmann, P. Kambadur, D. Rosenberg, and G. Mann, "Bloomberggpt: A large language model for finance," 2023.
- [46] P. Malo, A. Sinha, P. Korhonen, J. Wallenius, and P. Takala, "Good debt or bad debt: Detecting semantic orientations in economic texts," *Journal of the Association for Information Science and Technology*, vol. 65, 2014.
- [47] D. Roozen and F. Lelli, "Stock values and earnings call transcripts: a sentiment analysis," *Preprints 2021*, 2021020424, 2021. [Online]. Available: <https://dataverse.nl/dataset.xhtml?persistentId=doi:10.34894/TJE0DO>
- [48] A. Sinha and T. Khandait, "Impact of news on the commodity market: Dataset and results," 2020.
- [49] J. C. Salinas Alvarado, K. Verspoor, and T. Baldwin, "Domain adaption of named entity recognition to support credit risk assessment," in *Proceedings of the Australasian Language Technology Association Workshop 2015*, Parramatta, Australia, Dec. 2015, pp. 84–90. [Online]. Available: <https://aclanthology.org/U15-1010>
- [50] T. Deußer, S. M. Ali, L. Hillebrand, D. Nurchalifah, B. Jacob, C. Bauckhage, and R. Sifa, "KPI-EDGAR: A novel dataset and accompanying metric for relation extraction from financial documents," in *2022 21st IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, dec 2022. [Online]. Available: <https://doi.org/10.1109%2Ficmla55696.2022.00254>
- [51] L. Loukas, M. Fergadiotis, I. Chalkidis, E. Spyropoulou, P. Malakasiotis, I. Androulopoulos, and P. George, "Finer: Financial numeric entity recognition for xbrl tagging," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL 2022)*. Association for Computational Linguistics, 2022. [Online]. Available: <https://arxiv.org/abs/2203.06482>
- [52] [Https://sites.google.com/view/fiqqa/home](https://sites.google.com/view/fiqqa/home).
- [53] M. Feng, B. Xiang, M. R. Glass, L. Wang, and B. Zhou, "Applying deep learning to answer selection: A study and an open task," 2015.
- [54] Z. Chen, S. Li, C. Smiley, Z. Ma, S. Shah, and W. Y. Wang, "Convfinqa: Exploring the chain of numerical reasoning in conversational finance question answering," 2022.
- [55] Z. Zhou, L. Ma, and H. Liu, "Trade the event: Corporate events detection for news-based event-driven trading," 2021.
- [56] S. Biderman, H. Schoelkopf, Q. Anthony, H. Bradley, K. O'Brien, E. Hallahan, M. A. Khan, S. Purohit, U. S. Prashanth, E. Raff, A. Skowron, L. Sutawika, and O. van der Wal, "Pythia: A suite for analyzing large language models across training and scaling," 2023.

- [57] Y. Tay, M. Dehghani, V. Q. Tran, X. Garcia, J. Wei, X. Wang, H. W. Chung, S. Shakeri, D. Bahri, T. Schuster, H. S. Zheng, D. Zhou, N. Houlsby, and D. Metzler, “Ul2: Unifying language learning paradigms,” 2023.
- [58] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, “LLaMA: Open and efficient foundation language models,” *arXiv:2302.13971*, 2023.
- [59] L. Weidinger, J. Uesato, M. Rauh, C. Griffin, P.-S. Huang, J. Mellor, A. Glaese, M. Cheng, B. Balle, A. Kasirzadeh, C. Biles, S. Brown, Z. Kenton, W. Hawkins, T. Stepleton, A. Birhane, L. A. Hendricks, L. Rimell, W. Isaac, J. Haas, S. Legassick, G. Irving, and I. Gabriel, “Taxonomy of risks posed by language models,” in *Proceedings of the ACM Conference on Fairness, Accountability, and Transparency*, 2022, pp. 214–229.
- [60] R. Shelby, S. Rismani, K. Henne, A. Moon, N. Rostamzadeh, P. Nicholas, N. Yilla, J. Gallegos, A. Smart, E. Garcia, and G. Virk, “Sociotechnical harms: Scoping a taxonomy for harm reduction,” *arXiv preprint arXiv:2210.05791*, 2022.
- [61] IBM Corporation. IBM AI Ethics. [Online]. Available: <https://www.ibm.com/impact/ai-ethics>
- [62] B. Green, D. Heider, K. Firth-Butterfield, and D. Lim, “Responsible use of technology: The IBM case study,” World Economic Forum, White Paper, Sep. 2021.
- [63] “Foundation models: Opportunities, risks and mitigations,” IBM AI Ethics Board, Tech. Rep., Jul. 2023.
- [64] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, “On the dangers of stochastic parrots: Can language models be too big?” in *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, 2021, pp. 610–623.
- [65] L. Lessig, “The new chicago school,” *The Journal of Legal Studies*, vol. 27, no. S2, pp. 661–691, 1998.
- [66] H. R. Kirk, B. Vidgen, P. Röttger, and S. A. Hale, “Personalisation within bounds: A risk taxonomy and policy framework for the alignment of large language models with personalised feedback,” *arXiv preprint arXiv:2303.05453*, 2023.
- [67] I. D. Raji, E. Denton, E. M. Bender, A. Hanna, and A. Paullada, “AI and the everything in the whole wide world benchmark,” in *Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021.
- [68] S. Fazelpour and M. De-Arteaga, “Diversity in sociotechnical machine learning systems,” *Big Data & Society*, vol. 9, no. 1, p. 20539517221082027, 2022.
- [69] IBM Corporation. IBM Natural Conversation Framework. [Online]. Available: <https://ibm.biz/natconv>
- [70] R. J. Moore, S. An, and G.-J. Ren, “The IBM natural conversation framework: a new paradigm for conversational UX design,” *Human Computer Interaction*, vol. 38, no. 3-4, pp. 168–193, 2023. [Online]. Available: <https://doi.org/10.1080/07370024.2022.2081571>
- [71] M. Amaral, H. Chen, T. Chiba, R. Nakazawa, S. Choochotkaew, E. K. Lee, and T. Eilam, “Kepler: A framework to calculate the energy consumption of containerized applications,” in *IEEE International Conference on Cloud Computing*, 2023.
- [72] P. Sattigeri, S. Ghosh, I. Padhi, P. Dognin, and K. R. Varshney, “Fair infinitesimal jackknife: Mitigating the influence of biased training data points without refitting,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 35 894–35 906, 2022.
- [73] G. Zhang, Y. Zhang, Y. Zhang, W. Fan, Q. Li, S. Liu, and S. Chang, “Fairness reprogramming,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 34 347–34 362, 2022.
- [74] S. Basu, P. Sattigeri, K. N. Ramamurthy, V. Chenthamarakshan, K. R. Varshney, L. R. Varshney, and P. Das, “Equi-tuning: Group equivariant fine-tuning of pretrained models,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 6, 2023, pp. 6788–6796.