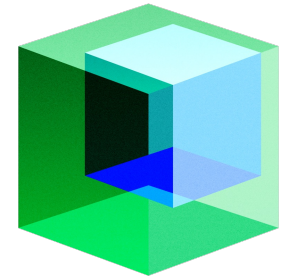


Open-Source Data & Tools for the IBM Granite AI Models

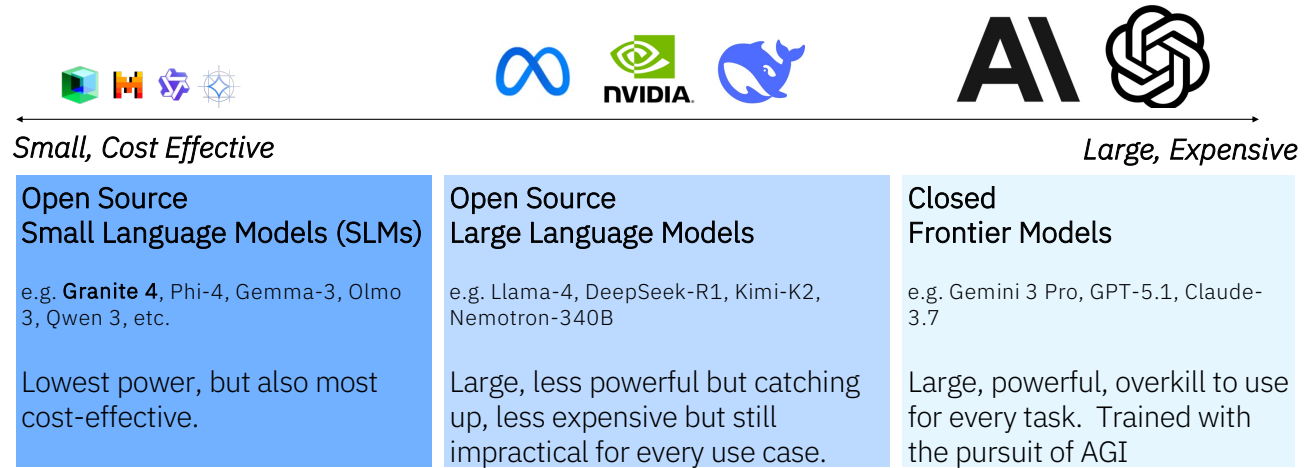
Petros Zerfos, Ph.D.

pzerfos@us.ibm.com

Principal Research Scientist & Manager
Data & Tools for AI Models
IBM Research

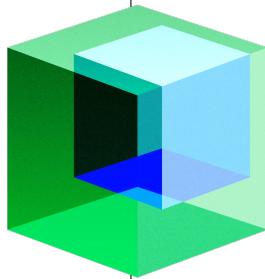


The Generative AI landscape is rapidly evolving



Granite

A family of enterprise-focused SLMs designed to help companies put GenAI to work.



Language

Granite-3.3-8B-Instruct
Granite-3.3-2B-Instruct

Granite 4.0 Family (NEW)



Safety Guardrails

Granite-Guardian-3.3-8B



Embeddings

Granite-Embeddings-30M-English
Granite-Embeddings-125M-English
Granite-Embeddings-107M-Multilingual
Granite-Embeddings-278M-Multilingual
Granite-Vision-3.3-2B-Embedding



Time Series

Granite-TimeSeries-Flowstate-r1.0



Vision

Granite-Vision-3.3-2B

Granite-Docting-258M (NEW)

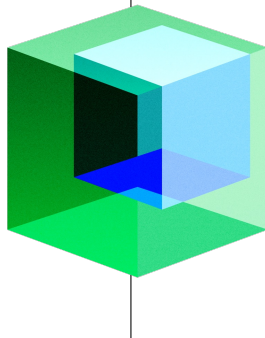


Speech

Granite-Speech-3.3-8B
Granite-Speech-3.3-2B

Granite 4

A family of enterprise-focused [SLMs](#) designed to help companies put generative AI to work.



© 2025 IBM Corporation



Open

- Everything Granite available under the no-nonsense [Apache 2.0](#) license
- Both instruct and base models shared for [easy customization](#)



Efficient

- Hybrid Mamba2 architecture designed for [small GPU footprint](#) and hardware-constrained deployments
- Range of model sizes provided to enable [optimizing model size for a use case](#)



Trusted

- First Open-Source model that is [ISO 42001 certified](#)
- Checkpoints [cryptographically-signed](#) to verify source of origin
- Launched with new white-hat hacker [bug bounty program](#)



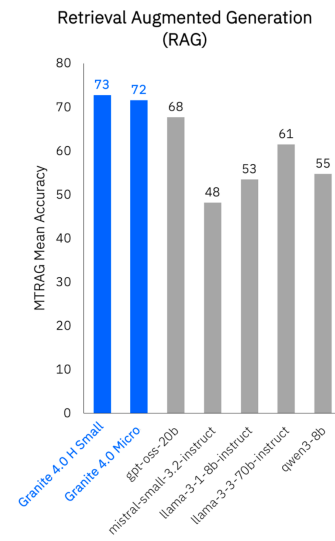
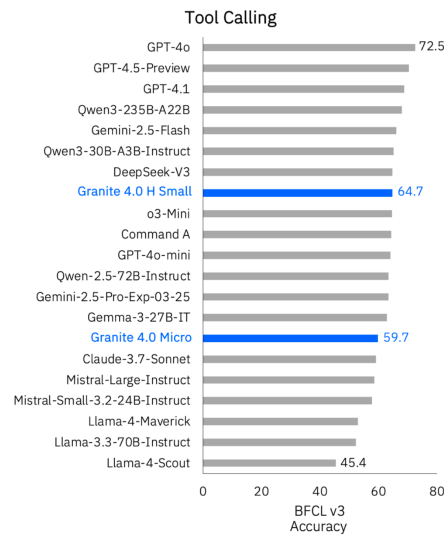
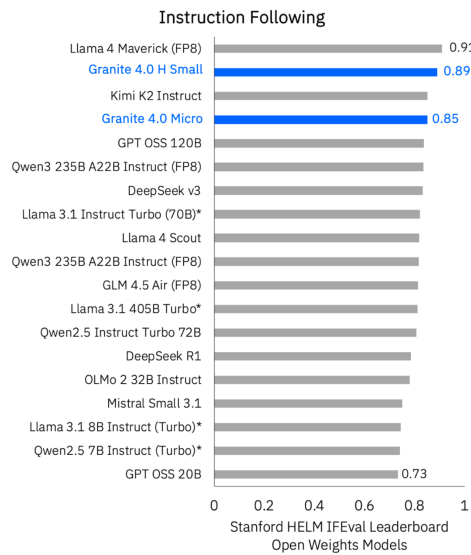
Granite 4

A family of enterprise-focused SLMs designed to help companies put generative AI to work.



	Granite-4.0-H-Small	Granite-4.0-H-Tiny	Granite-4.0-H-Micro	Granite-4.0-Micro
Architecture Type	Hybrid, Mixture of Experts	Hybrid, Mixture of Experts	Hybrid, Dense	Traditional, Dense
Model Size	32B total parameters 9B activated parameters	7B total parameters 1B activated parameters	3B total parameters	3B total parameters
Intended Use	Workhorse model for key enterprise tasks like RAG and agents	Designed for low latency, edge, and local applications, and as a building block to perform key tasks (like function calling) quickly within agentic workflows*		Alternative option for users when Mamba2 support is not yet optimized (e.g. llama.cpp, PEFT, etc)
Est Memory Reqs (8-bit, 128K context length, batch = 1)	33 GB	8 GB	4 GB	9 GB
Example Hardware** (8-bit, 128K context length, batch = 1)	NVIDIA L40S (<\$10K)	RTX 3060 12GB (<\$1K)	Raspberry-Pi 8GB (<\$100)	RTX 3060 12GB (<\$1K)

Granite 4.0 Optimizes Performance for Enterprise-Relevant Tasks



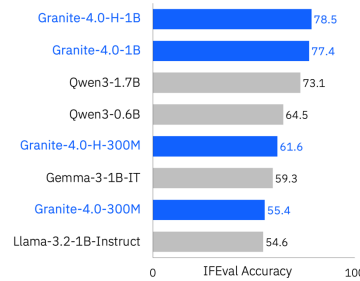
© 2025 IBM Corporation

*Original model is open weights, but exact Turbo version could not be found on HF

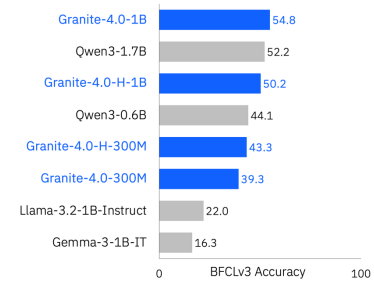
Granite 4.0 Nano

- **Granite 4.0 H 1B** – A ~1.5B parameter, dense LLM featuring a hybrid-SSM based architecture.
- **Granite 4.0 H 350M** – A ~350M parameter, dense LLM featuring a hybrid-SSM based architecture.
- **Granite 4.0 1B and Granite 4.0 350M** – Alternative traditional transformer versions of our 1B and 350M Nano models

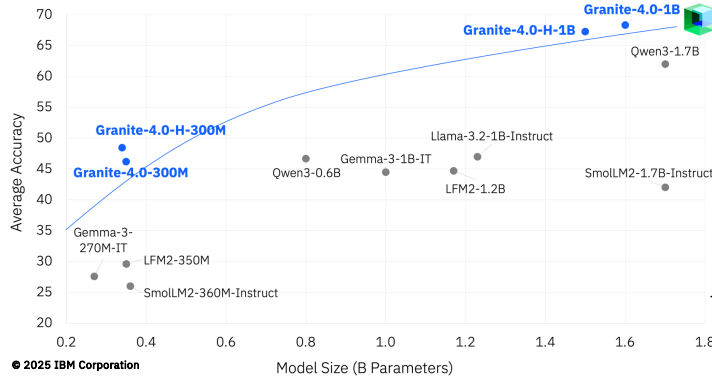
Instruction Following



Tool Calling

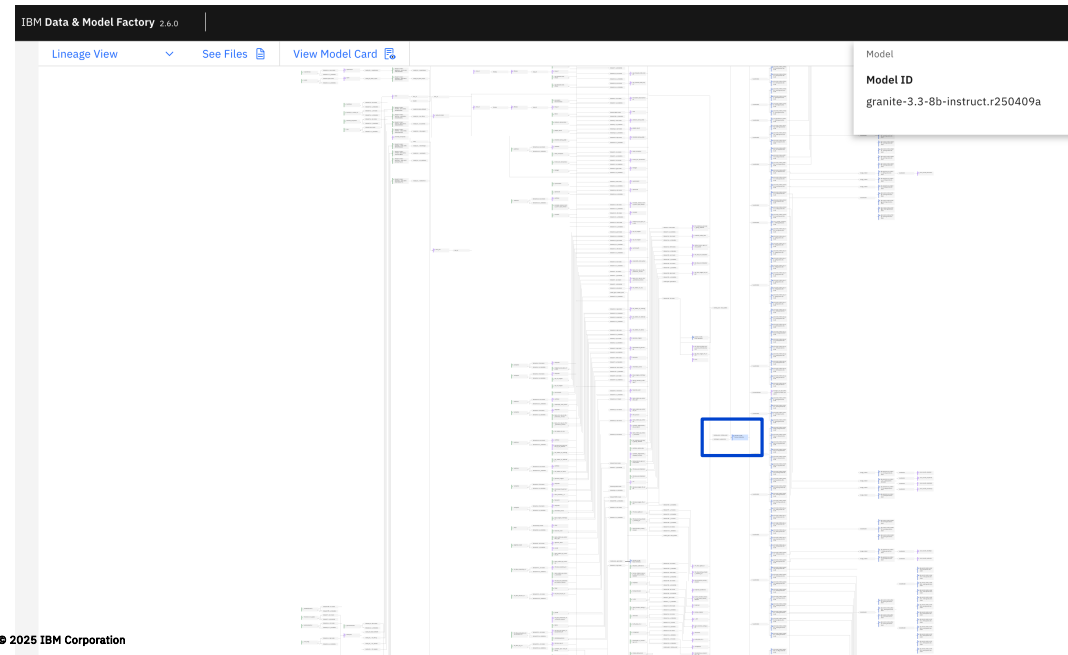


General Performance vs Model Size



Benchmarks across **General Knowledge** (MMLU, BBH), **Math** (GSM8K, GSM8K-Symbolic), **Code** (EvalPlus, CruxEval-O) and **Safety** (AttaQ, SALAD-Bench)

Trust depends on data and model governance





IBM doubles down on trust & safety partnerships for Granite



Ranked 1st in Stanford AI Transparency Index 2025 (from #4 in 2024)

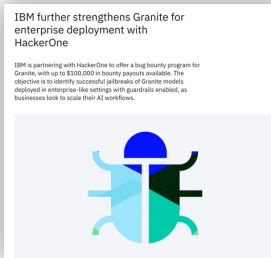
New \$100k bug bounty program in partnership with HackerOne

New red teaming partnership for Granite-vision with HiddenLayer

New partnership with Schellman LLC. Granite obtained ISO 42001 compliance certification

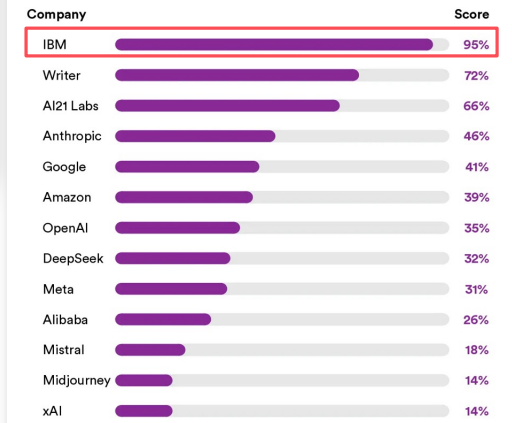
Renewed red teaming partnership with Robust Intelligence

© 2025 IBM Corporation



Foundation Model Transparency Index Total Scores, 2025

Source: 2025 Foundation Model Transparency Index



Stanford AI Transparency Index 2025

<https://crfm.stanford.edu/fmti/December-2025/index.html>

Every conversation in AI starts with Models and ends with Data

79% identify data preparation and generation* as the most common strategic task performed by AI teams.

30% view data volume and complexity* as one of the most challenging aspects of AI implementation.

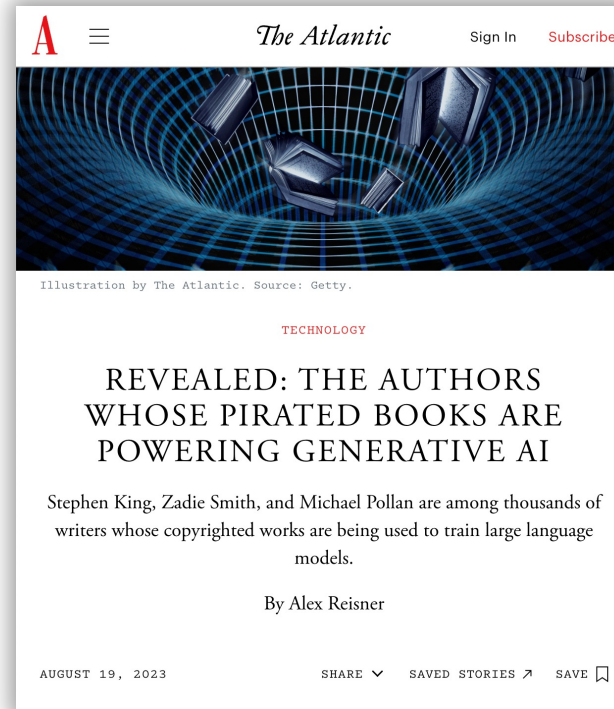
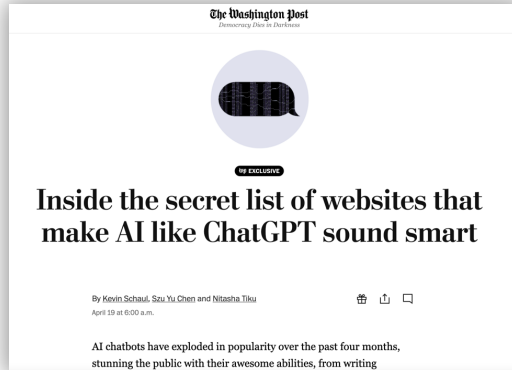
Quality of data affects the quality of the model


* Gartner, Explore Data-Centric AI Solutions to Streamline AI Development, 2023



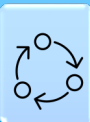
George Fuechsel,
IBM 305 RAMAC
technician

Issues with GenAI training data curation



 The Data:
GneissWeb

 The Tool:
Data Prep Kit

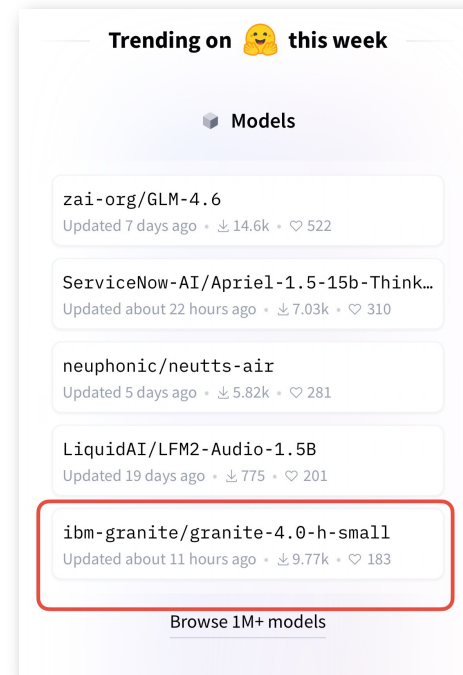
 The Ops:
IBM Cloud

GneissWeb was used to train IBM Granite 4.0

- Trending in Top-5 at HuggingFace 🤗
(as of Oct. 7th)
- Granite 4.0:
 - Granite-4.0-H-Small (MoE 32B/A9B)
 - Granite-4.0-H-Tiny (MoE 7B/A1B)
 - Granite-4.0-H-Micro (Dense 3B)

IBM Granite 4.0 Release:

<https://www.ibm.com/new/announcements/ibm-granite-4-0-hyper-efficient-high-performance-hybrid-models>

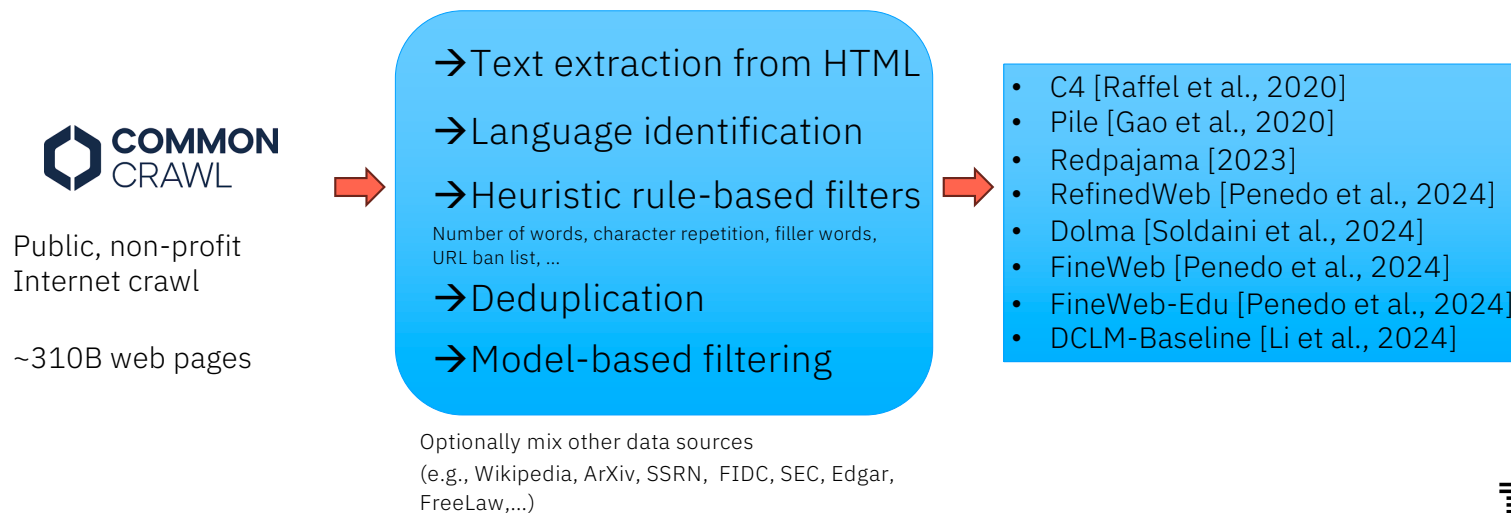


The screenshot shows the 'Trending on this week' section of the HuggingFace models page. It lists five models under the 'Models' category. The model 'ibm-granite/granite-4.0-h-small' is highlighted with a red box. Below the list is a link to 'Browse 1M+ models'.

Model Name	Updated	Downloads	Likes
zai-org/GLM-4.6	Updated 7 days ago	14.6k	522
ServiceNow-AI/Apriel-1.5-15b-Think...	Updated about 22 hours ago	7.03k	310
neuphonic/neutts-air	Updated 5 days ago	5.82k	281
LiquidAI/LFM2-Audio-1.5B	Updated 19 days ago	775	201
ibm-granite/granite-4.0-h-small	Updated about 11 hours ago	9.77k	183

Good models need Good data!

Open datasets for pre-training LLMs

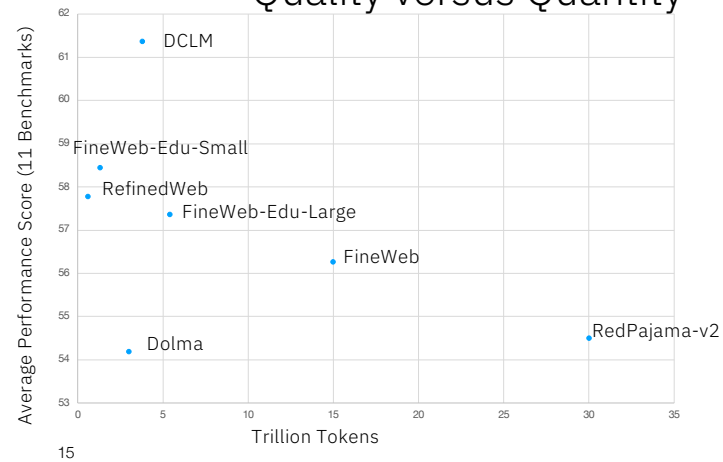


Limitations of current pre-training datasets

State-of-the-Art Datasets

- Large datasets (5T+ tokens) – limited performance
- Aggressive filtering to achieve quality
- Rely on model-based filtering for quality, often requiring GPUs

Quality versus Quantity



GneissWeb: advancing open innovation in large-scale training data

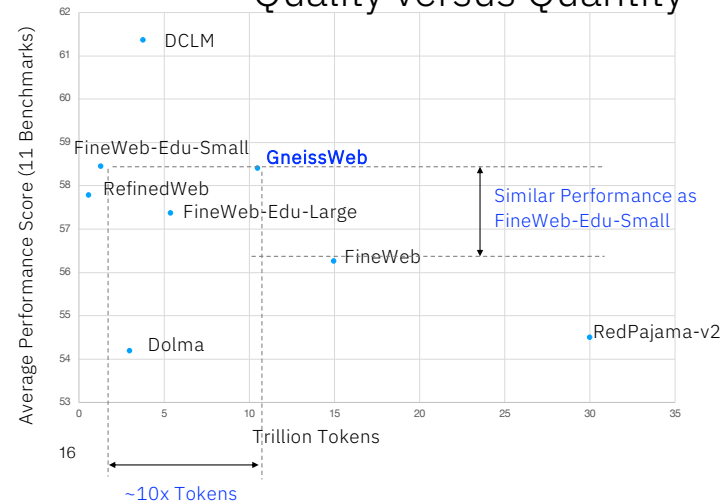
State-of-the-Art Datasets

- Large datasets (5T+ tokens) – limited performance
- Aggressive filtering to achieve quality
- Rely on model-based filtering for quality, often requiring GPUs

GneissWeb

- 10T tokens / 12.6B docs / 35TB size with quality higher than large (5T+ tokens) SOTA datasets
- CPU-friendly recipes (model+heuristics)
- Novel quality annotators, judiciously designed ensemble filter
- Finance data: ~5TB

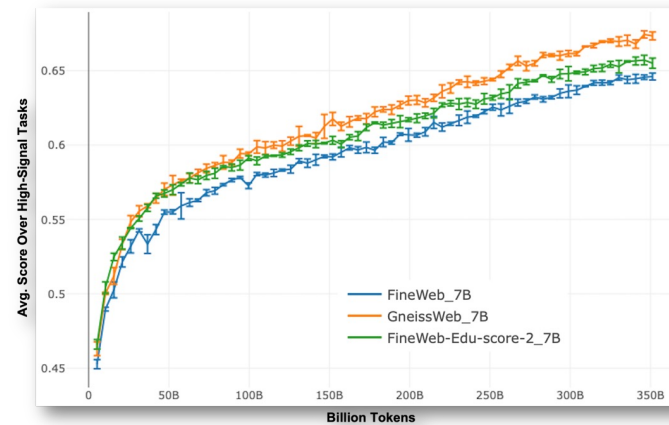
Quality versus Quantity



GneissWeb: SOTA Common Crawl for pre-training

- 10T token dataset derived from FineWeb V1.1.0
- Ablation models *outperform* those trained on FineWeb V1.1.0 by **>1.5%** over 20 benchmarks
- Open-source tools and recipes for reproduction

<https://huggingface.co/datasets/ibm-granite/GneissWeb>



Combining GneissWeb components into a winning recipe

Exact Substring Deduplication

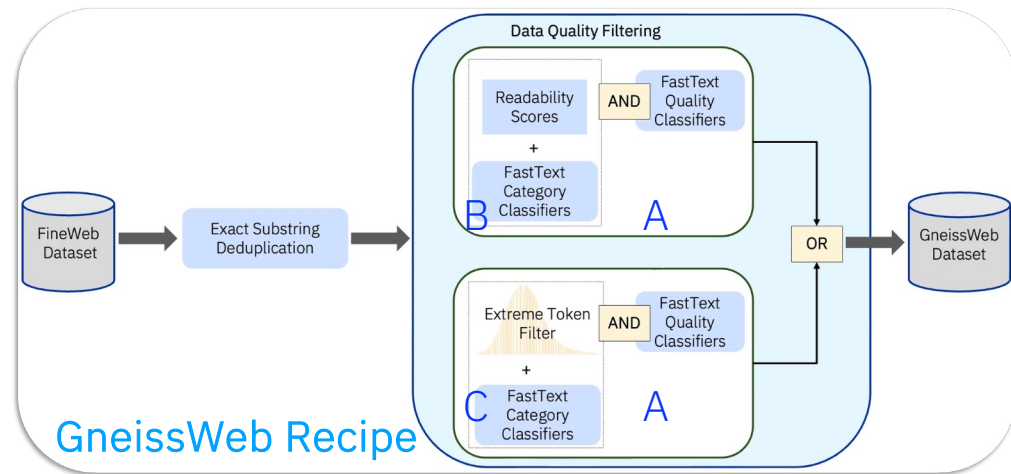
Custom Data Quality Classifiers

Filtering based on Readability Scores

Filtering Extreme-Tokenized Documents

Document Category Classifiers

Novel Components highlighted in ■

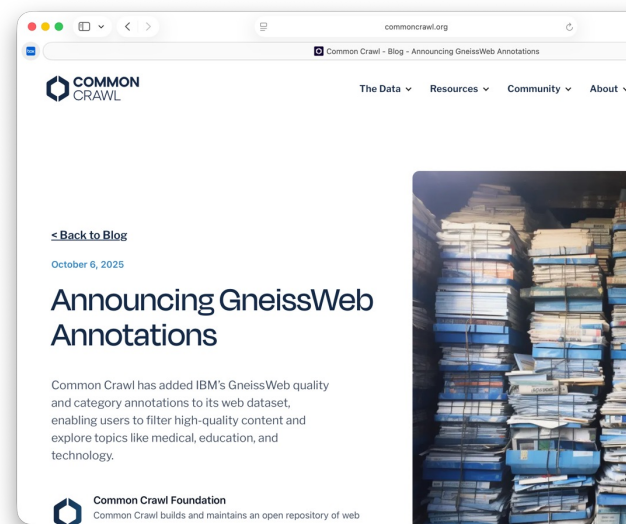


GneissWeb Recipe

- Exact substring deduplication → ((A AND B) OR (A AND C))
- GneissWeb ensemble filtering: A document is retained if either the fastText combination and category-aware, readability score filter agree to retain,
- OR the fastText combination and category-aware, extreme-tokenized filter agree to retain

Common Crawl Foundation is bringing GneissWeb to everyone!

- Announced at IBM TechXchange 2025
- Using IBM's Bloom filter + Data Prep Kit (LF AI & Data Project)
- Applied it to every URL in the Crawls
- Created annotations
 - Quality score: passes GneissWeb's standard or not
 - Category: science, medical, education, etc.
- Published annotations at both URL- and host-level

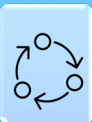


Common Crawl Foundation Announcement:

<https://commoncrawl.org/blog/announcing-gneissweb-annotations>

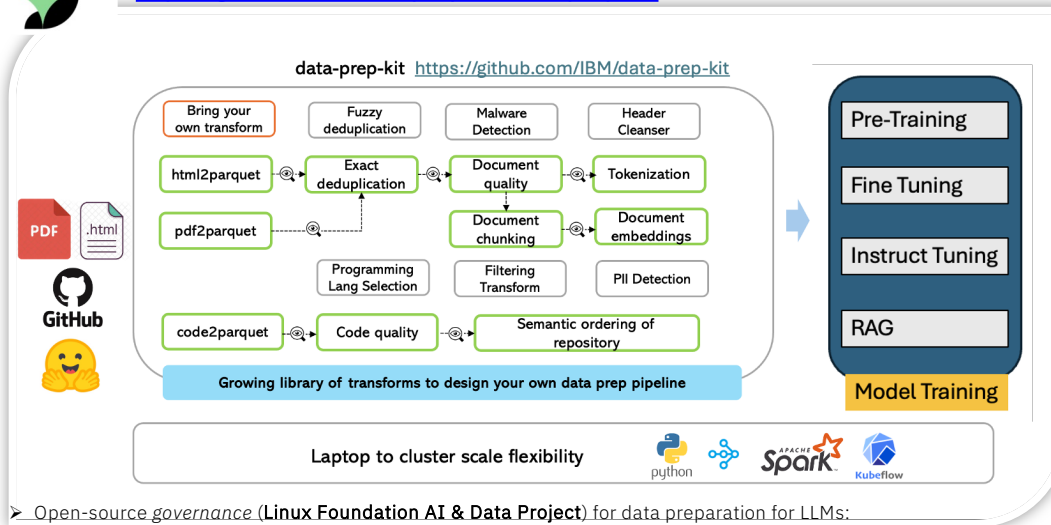
 **The Data:**
GneissWeb

 **The Tool:**
Data Prep Kit

 **The Ops:**
IBM Cloud

Data Prep Kit

<https://github.com/data-prep-kit/data-prep-kit>



- Open-source *governance* (Linux Foundation AI & Data Project) for data preparation for LLMs: <https://lfaidata.foundation/projects/data-prep-kit/>
- Spearheaded by IBM Research
- Used in the data preparation for IBM's Granite AI model training: <https://huggingface.co/ibm-granite>

Large
collection of
transforms

**40+ built in
Transforms**

Bring Your Own Transform



Content Extraction (including Granite-Docling)

Document and Code quality

Data Enrichment

Annotations and filtering

Language Identification

Quality Annotation

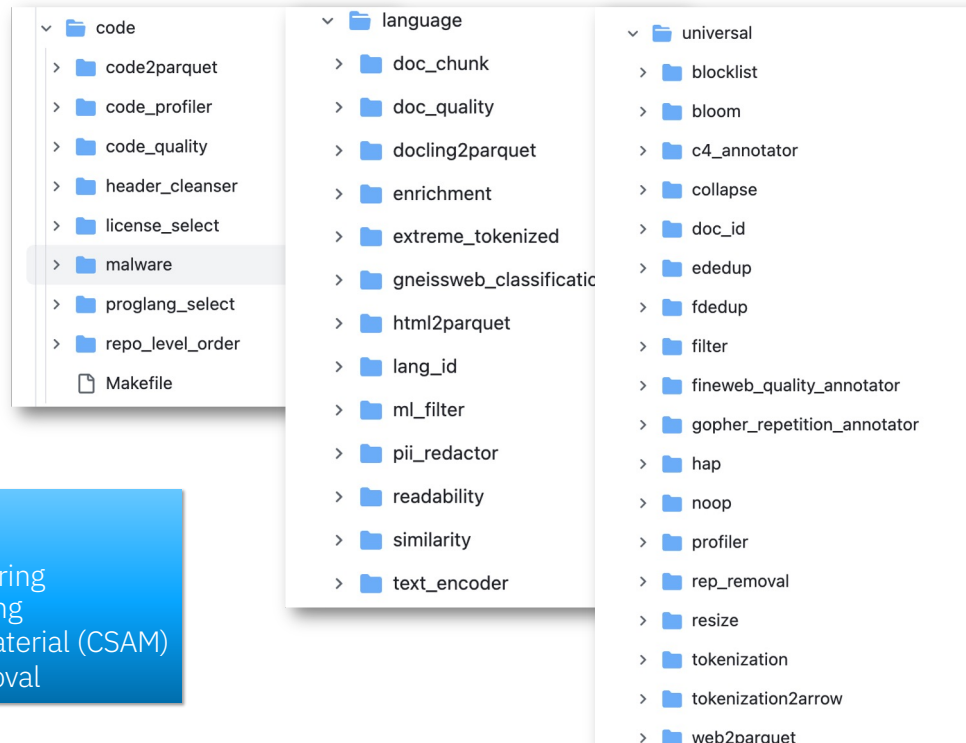
Deduplication

Data Safety and Security

Optimization & ready for training

Bring your own Transform

Data Prep Kit Transforms










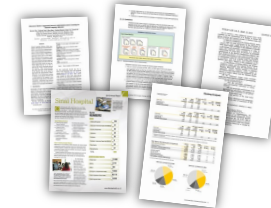
Granite Vision (VLM)

- Face blurring
- License-based filtering
- Non-US data filtering
- Child Sex Abuse Material (CSAM) detection and removal

Introducing Docling

- 1. **Open source** and permissively licensed (**Apache 2.0**)
- 2. **Cost performant**
- 3. Packaged as a **python library** (**no API's**)

-  Parsing of multiple document formats incl. PDF, DOCX, XLSX, HTML, images, and more
-  Unified, expressive DoclingDocument representation format
-  Various export formats (Markdown, HTML, JSON)
-  Many plug-and-play ecosystem integrations
-  Support of OCR and Visual Language Models
-  Support for Audio with Automatic Speech Recognition models
-  Simple and convenient CLI



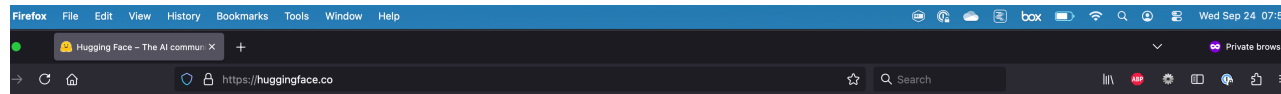
```
pip install docling

# a single document to markdown
docling https://arxiv.org/pdf/2408.09869.pdf

# a folder of documents to markdown and json
docling --to json --to md ./inputs/
```



Granite-Docling Adoption



Trending on 🤖 this week

Models	Spaces	Datasets
<p>ibm-granite/granite-docling-258M Updated about 21 hours ago · ⬆️ 25.7k · ❤️ 593</p>	<p>DeepSite v2 🌟 Generate any application with DeepSeek · ❤️ 142k</p>	<p>HuggingFaceFW/finepdfs Updated 16 days ago · ⬆️ 80.2k · ❤️ 564</p>
<p>openbmb/VoxCPM-0.5B Updated 5 days ago · ⬆️ 3.24k · ❤️ 649</p>	<p>Wan2.2 Animate 🌟 Wan2.2 Animate · ❤️ 419</p>	<p>fka/awesome-chatgpt-prompts Updated Jan 6 · ⬆️ 48.4k · ❤️ 9.12k</p>
<p>Alibaba-NLP/Tongyi-DeepResearch-30B-A3B Updated 7 days ago · ⬆️ 9.15k · ❤️ 601</p>	<p>Wan2.2 14B Fast 🌟 generate a video from an image with a text prompt · ❤️ 1.25k</p>	<p>InternRobotics/OmniWorld Updated 26 minutes ago · ⬆️ 15.2k · ❤️ 59</p>
<p>Wan-AI/Wan2.2-Animate-14B Updated 5 days ago · ⬆️ 14.4k · ❤️ 379</p>	<p>granite-docling-258M demo 🌟 Convert images to structured documents and answer questions · ❤️ 158</p>	<p>LucasFang/FLUX-Reason-6M Updated 12 days ago · ⬆️ 38.5k · ❤️ 78</p>
<p>Qwen/Qwen3-Omni-30B-A3B-Instruct Updated 1 day ago · ⬆️ 1.69k · ❤️ 313</p>	<p>indexTTS 2 Demo 🌟 Generate expressive speech from text with emotion control · ❤️ 268</p>	<p>HuggingFaceM4/FineVision Updated 19 days ago · ⬆️ 257k · ❤️ 342</p>
Browse 1M+ models	Browse 400k+ applications	Browse 250k+ datasets

Preview for Doc Conversion

EXPERIMENTAL RESULTS OF WASP-121b FROM JWST/NIRISS PHASE CURVE

while the kernel weights are structured in $(N_{\text{obs}}, N_{\text{obs}})$. This precomputation significantly accelerates our calculations, which is essential since the likelihoods are at least partially dependent with one another. Consequently, the inference runs faster and we can increase the number of iterations.

In addition, we follow a similar approach to our standard fit using MCMC, but we increase the total number of steps to 100,000 and use 100 walkers. Naturally, the fit would be faster if we had a parameter N_{obs} for the albedo values, N_{obs} for the rotation parameters, and one additional free parameter, α . However, since right-to-left does not contribute to the reflected light component, we exclude these albedo values from the fit. In any case, our choice of 100 walkers ensures a sufficient number of walkers per free parameter. Following Calcutt et al. (2022) we set an upper prior limit of 1/2 on the albedo values and a lower prior limit of -1 on the rotation parameters. We assume a Gaussian prior for the thermal emission to improve a uniform prior between 0 and 100 days for the phase.

We choose to fit our detected lightcurve considering 4 and 8 longitudinal slices $(N_{\text{obs}} = 4, 8)$. However, we show the results of the simplest 4 slice model. As in our previous fits, we started an initial run with 25,000 steps (25% of the total run) and set the unconstrained parameters from the previous fit as the starting positions for the final 75,000 step run. We then discard the first 95% of the final run as burn-in.

2.3. Planetary Effective Temperature

Phase curves are the only way to probe thermal emission from the day and nightside of an exoplanet and have been extensively studied recently (e.g., Tregloan-Rep et al. 2023). The wavelength range of NIRISS-SCIENCE covers a large portion of the emitted flux of WASP-121b (~ 30 KHz), see Figure 3, enabling a precise and robust constraint of the planet's energy budget.

We assume the fitted $F_{\text{p},\lambda}$ emission spectra to brightness temperature by wavelength,

$$B_{\lambda, \text{planet}} = \frac{F_{\text{p},\lambda}}{\pi R_p^2 D^2} \cdot B_{\lambda, \text{star}} \quad (16)$$

where the planet's thermal emission is

$$B_{\lambda, \text{planet}} = \frac{F_{\text{p},\lambda}}{\pi R_p^2 D^2} \cdot B_{\lambda, \text{star}} \quad (17)$$

There are many ways of converting brightness temperature to effective temperature, including the Error-Weighted Mean (EWM), Power-Weighted Mean (PWM) and with a Gaussian Process (Gaussian & Cross 2015).

Finn et al. (2015). In this work, we elect to compute our effective temperature estimate with a novel method that is essentially a combination of the PWM and EWM. We create the effective temperature by using a Weighted Mean-Cut process. First, we generate our $F_{\text{p},\lambda}$ emission spectra at each point in the orbit by the Gaussian based on the measurement uncertainty. Our new emission spectrum is then used to create an estimate of the brightness temperature spectrum. The process is repeated at each orbital phase. We then estimate the effective temperature T_{eff} for a given orbital phase as

$$T_{\text{eff}} = \frac{\sum_{\lambda} w_{\lambda} B_{\lambda, \text{planet}}}{\sum_{\lambda} w_{\lambda}} \quad (18)$$

where w_{λ} is the weight for the slice wavelength given by the fraction of the planet's bolometric flux that falls within that wavelength bin, weighted by the inverse variance of the measurement,

$$w_{\lambda} = \frac{F_{\text{p},\lambda} \Delta \lambda \cdot \Delta \lambda \cdot \Delta \lambda}{\sum_{\lambda} F_{\text{p},\lambda} \Delta \lambda \cdot \Delta \lambda \cdot \Delta \lambda} \quad (19)$$

with T_{eff} representing an estimated effective temperature at the orbital phase of interest. When computing

Q&A about Doc Conversion

Convert this page to docling.

Does the document contain tables?

Can you extract the 2nd section header?

What element is located at <loc_84> <loc_403><loc_238><loc_419>

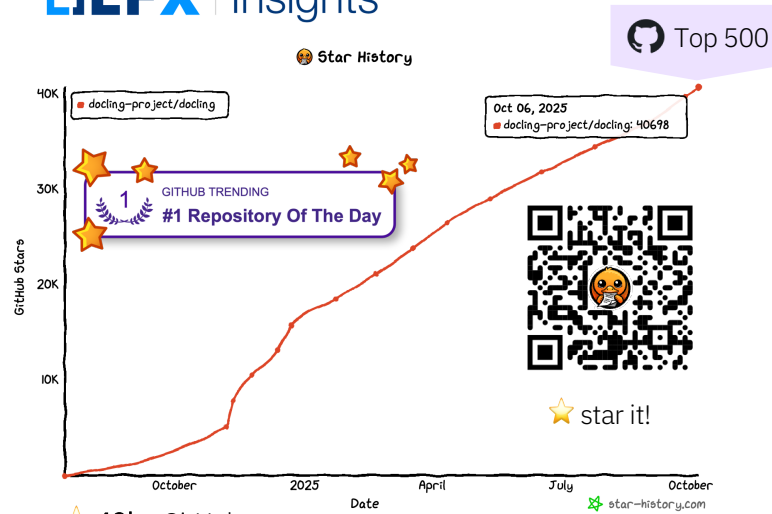
How can effective temperature be computed?

Extract all picture elements on the page.

Granite-Docling content extraction supports Arabic and right-to-left reading flow

Docling Community Adoption

OLFX | Insights

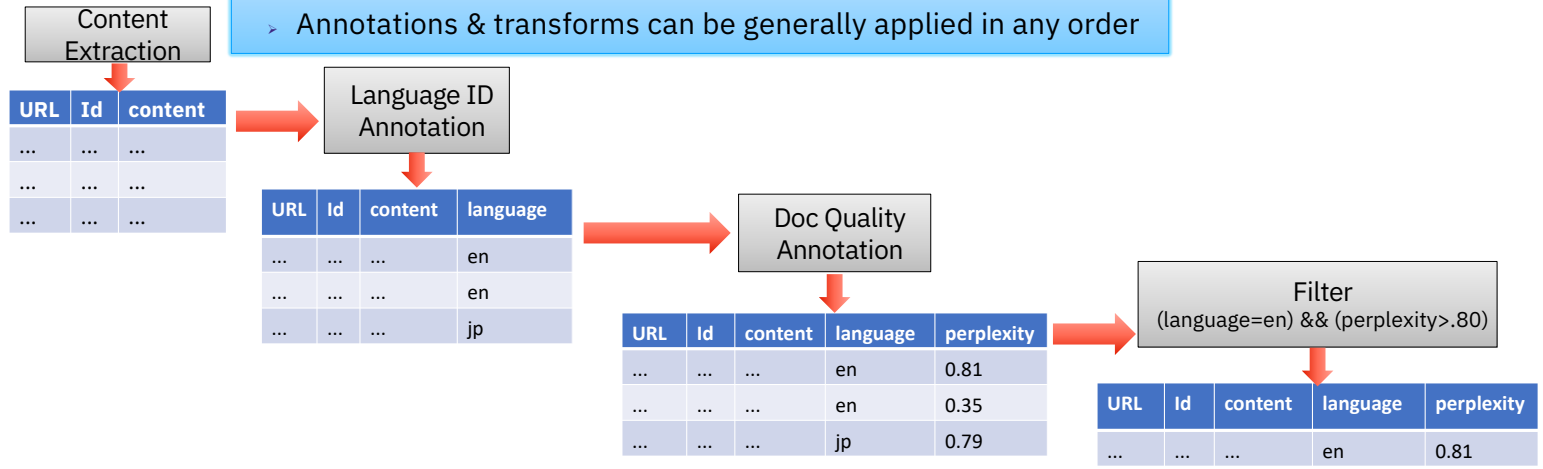


★ 42k+ GitHub stars

📦 1.5M+ downloads/month from [PyPi](#)

Unstructured Data Representation in Data Prep Kit

- Documents (language & code) are represented in a CSV-like form, in PyArrow tables (Parquet files)
- Annotations & transforms can be generally applied in any order



Data Prep Kit: from laptop- to cluster- scale



Laptop




Server



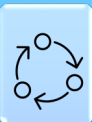
Datacenter

- ✓ Ray and Spark wrappers for scale-out with no code changes*
- ✓ Docker/Podman desktop, Kind cluster, local (file system) I/O, Cloud Object Storage, etc.
- ✓ Production-ready use: checkpointing, metadata, auditing
- ✓ Low-code pipeline orchestration via (KFP)

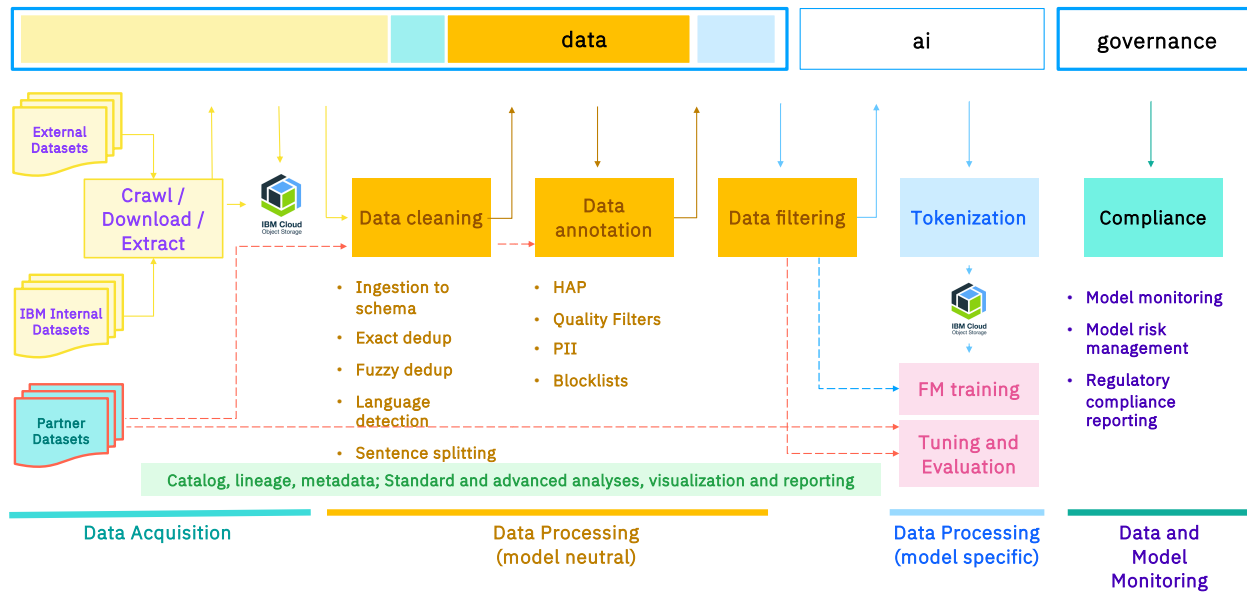
*for embarrassingly parallel transforms

 The Data:
GneissWeb

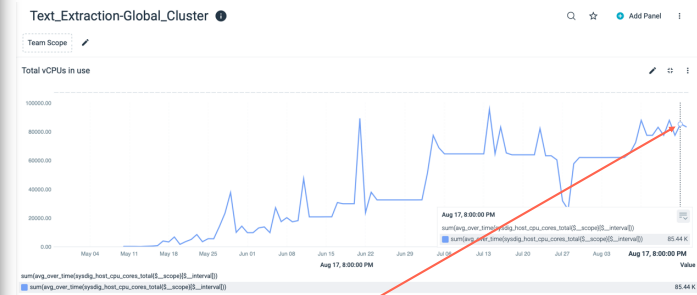
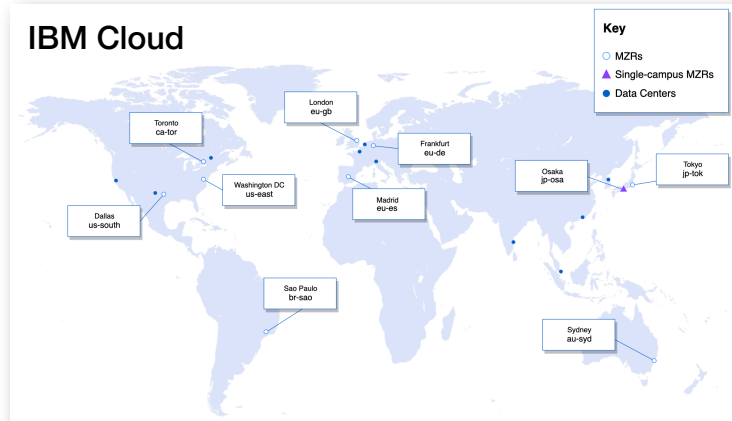
 The Tool:
Data Prep Kit

 The Ops:
IBM Cloud

Data engineering journey: from raw data to models



Data Engineering for LLMs @ IBM Cloud



Text Extraction on ~90K vCPUs across 2 regions / 6 data centers

Large-scale data processing with the IBM Cloud as “Supercomputer”

- text extraction from the HTML pages of the 12PBs of Common-Crawl dataset (~309B docs)

