# Unlocking AI Performance with NeMo Curator: Scalable Data Processing for LLMs

Ayush Maheshwari
Sr. Solutions Architect, NVIDIA

https://github.com/ayushbits/llm-development
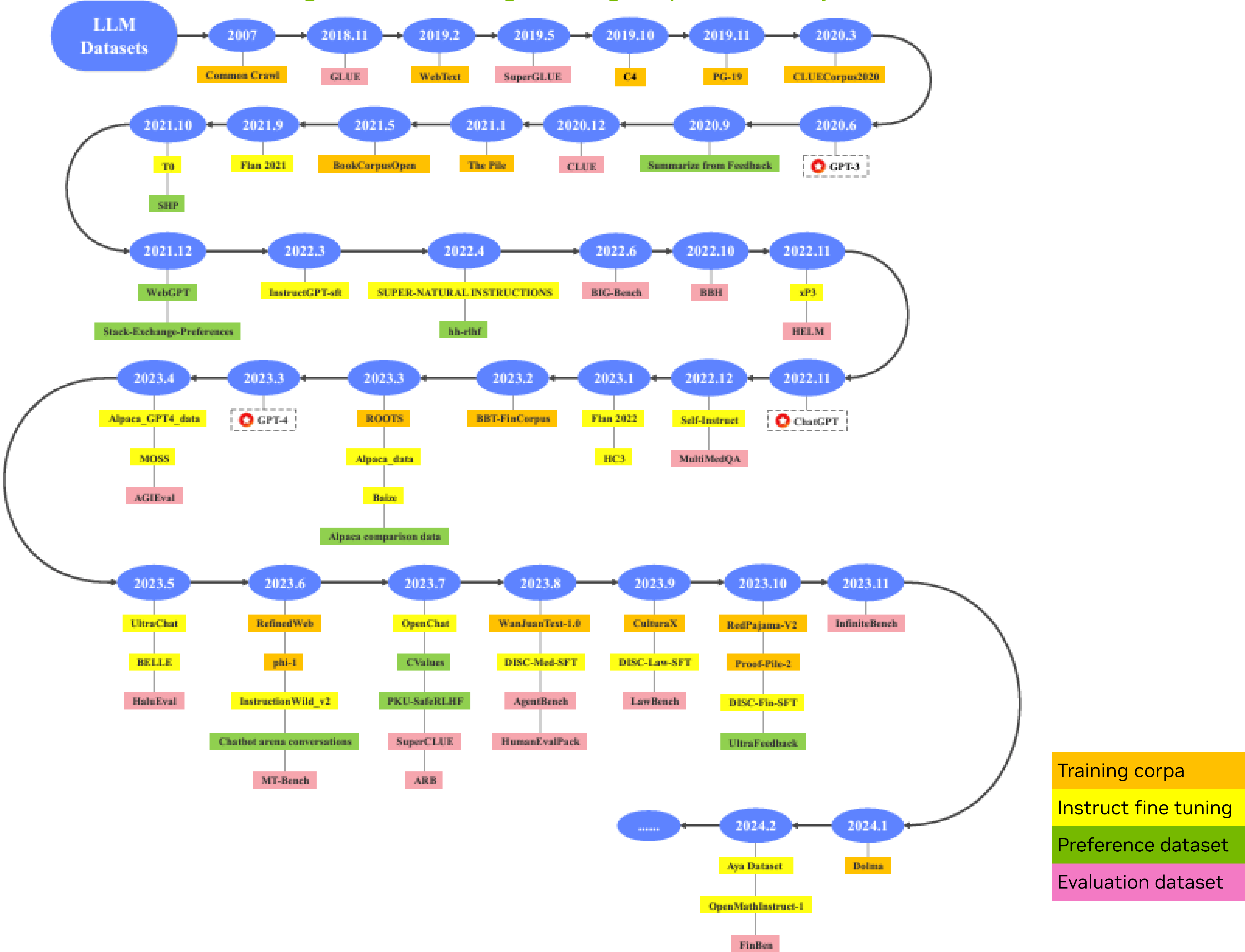
ayushbits.github.io

# Sessions

1. Cluster health-check using NCCL, MLPerf, HPL          **(1 hour) - Completed**
   a) Understand the hardware and its performance on multiple GPUs.
   b) Ensure that your training performance aligns with the h/w benchmarks
   c) Evaluate the cluster to ensure platform fits within your needs.

2. Large scale data curation for LLM training          **(1 hour) - Today**
   a) Deep-dive into aspects of data curation
   b) Mixed-precision training

3. Distributed and stable LLM training on a large-scale cluster     **(1.5 hour)**
   a) Parallelism techniques
   b) Frameworks and wrappers
   c) Recipes and best practices

4. Post-training and evaluation of pre-trained LLM          **(1 hour)**
   a) Sync between training data and expected performance
   b) Algorithms and frameworks

5. Fine-tuning and deployment          **(1 hour)**
   a) Dynamic and static batching, state management, inference server
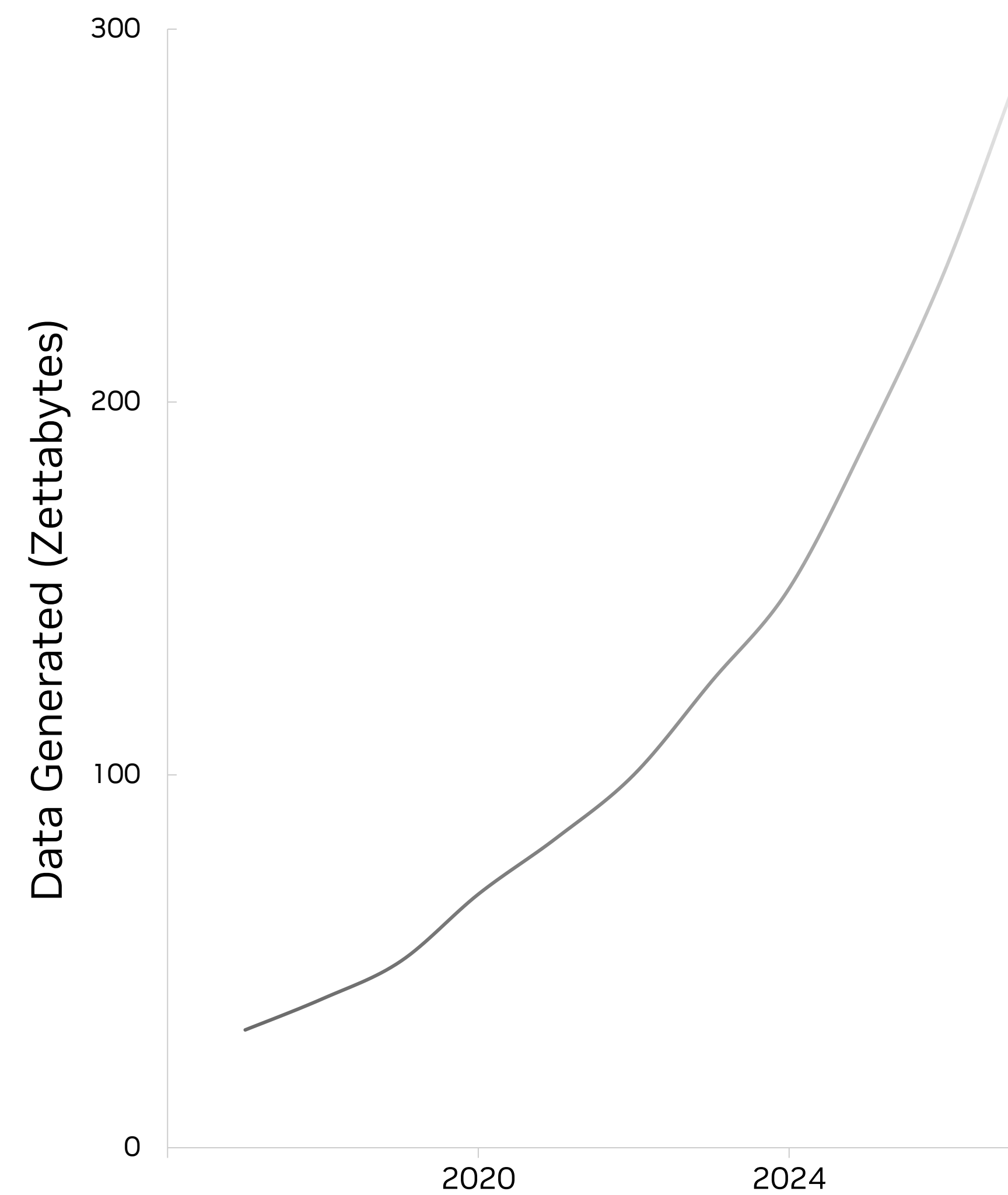   b) Best practices for optimizing model

NVIDIA.

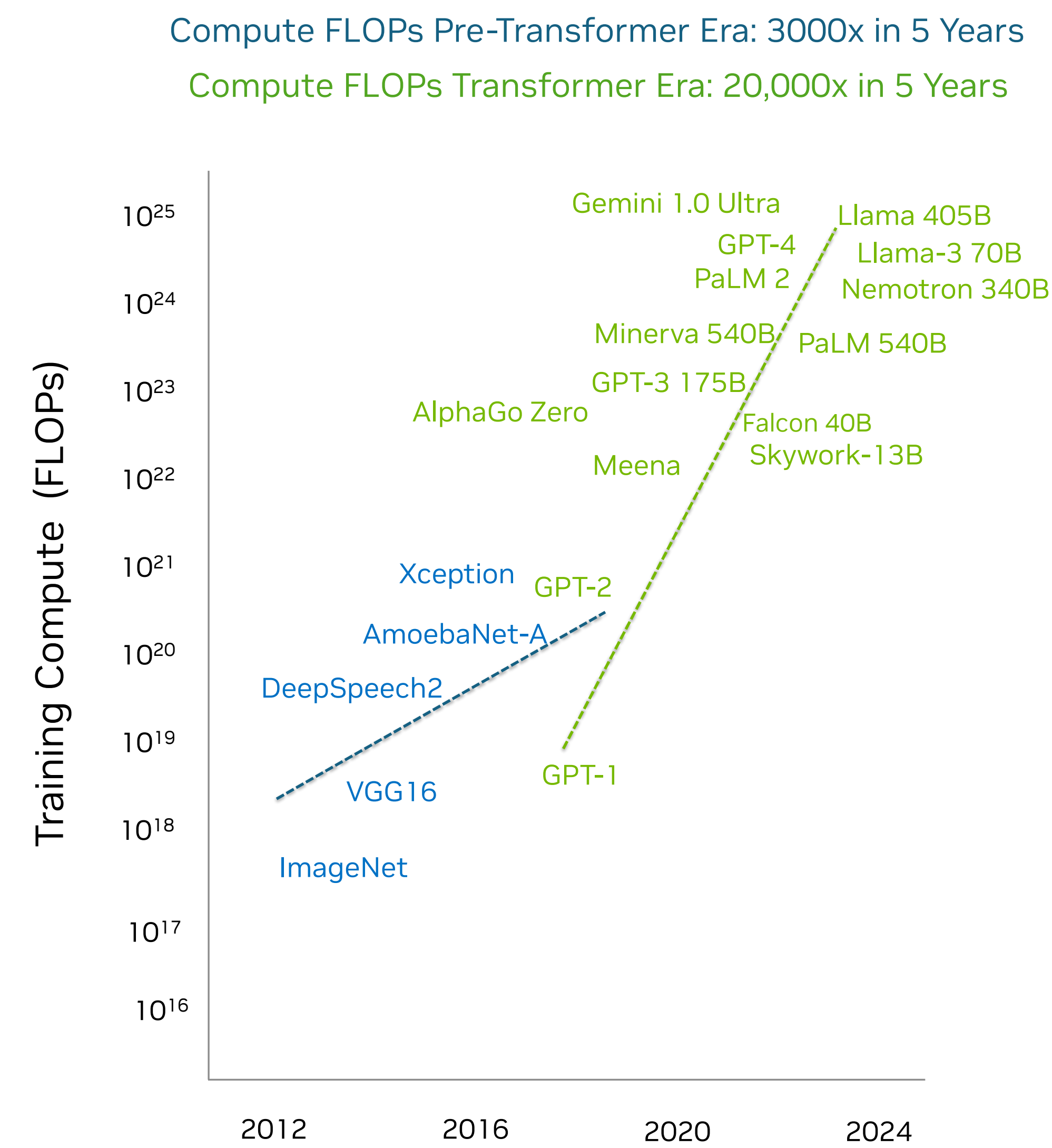# LLMs Are Trained on Internet Scale Data

## Data generated is growing exponentially
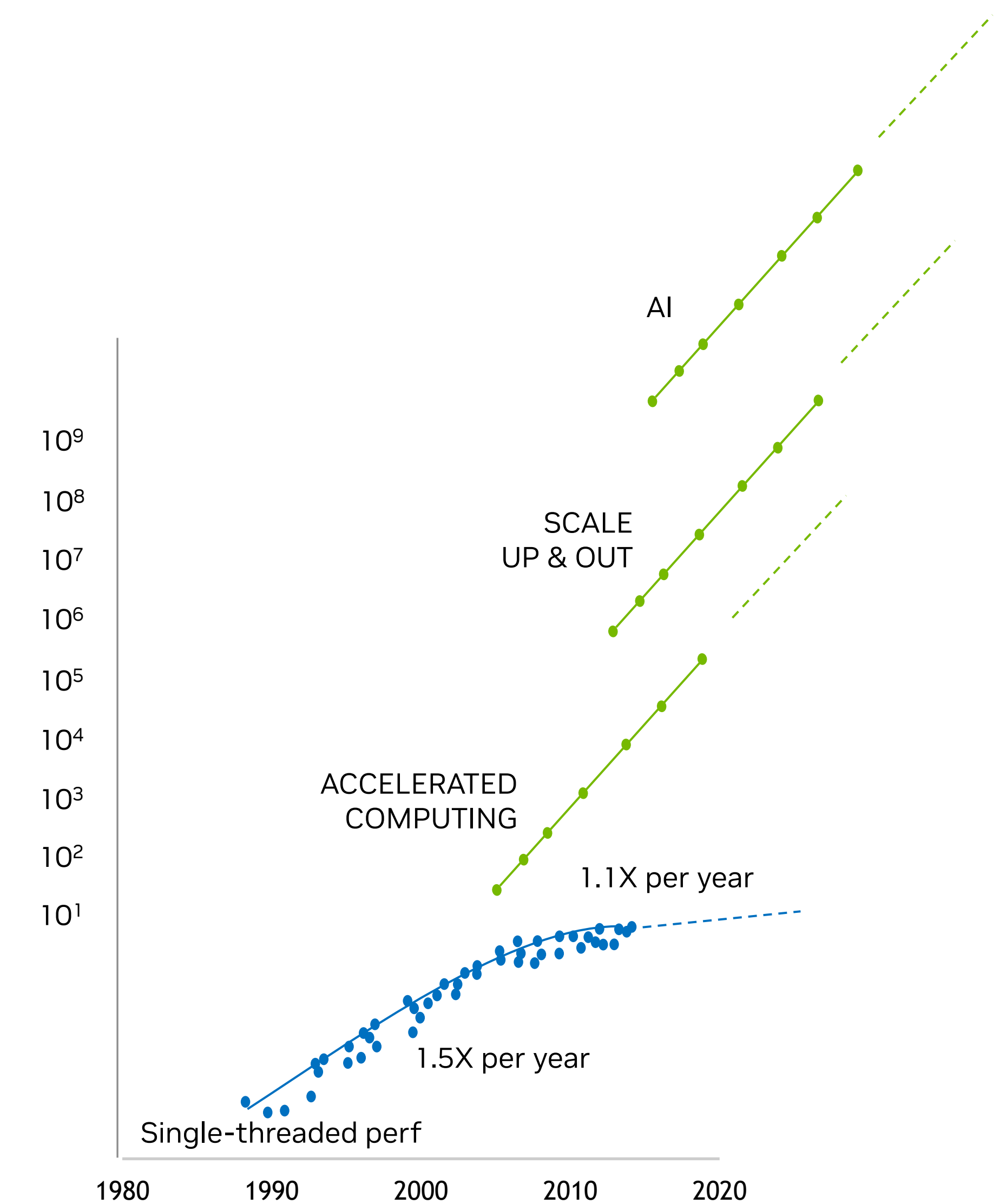
# Data Processing for LLMs Needs Accelerated Computing
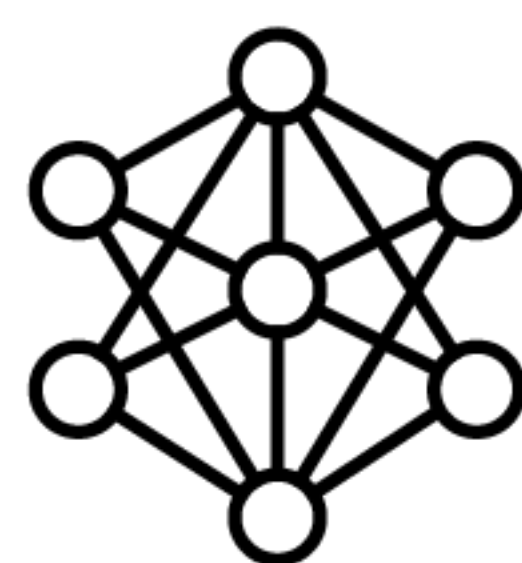
## Petabytes of Data Generated Yearly



Data Generated (Zettabytes)

300
200
100
0

2020    2024

## LLMs Trained on Internet Scale Data

Compute FLOPs Pre-Transformer Era: 3000x in 5 Years

Compute FLOPs Transformer Era: 20,000x in 5 Years



Training Compute  (FLOPs)

$10^{25}$
$10^{24}$
$10^{23}$
$10^{22}$
$10^{21}$
$10^{20}$
$10^{19}$
$10^{18}$
$10^{17}$
$10^{16}$

Gemini 1.0 Ultra    Llama 405B
GPT-4    Llama-3 70B
PaLM 2    Nemotron 340B
Minerva 540B    PaLM 540B
GPT-3 175B
AlphaGo Zero    Falcon 40B
Meena    Skywork-13B
Xception    GPT-2
AmoebaNet-A
DeepSpeech2
VGG16    GPT-1
ImageNet

2012    2016    2020    2024

## Moore's Law Has Ended



$10^9$
$10^8$
$10^7$
$10^6$
$10^5$
$10^4$
$10^3$
$10^2$
$10^1$

AI

SCALE
UP & OUT

ACCELERATED
COMPUTING

1.1X per year

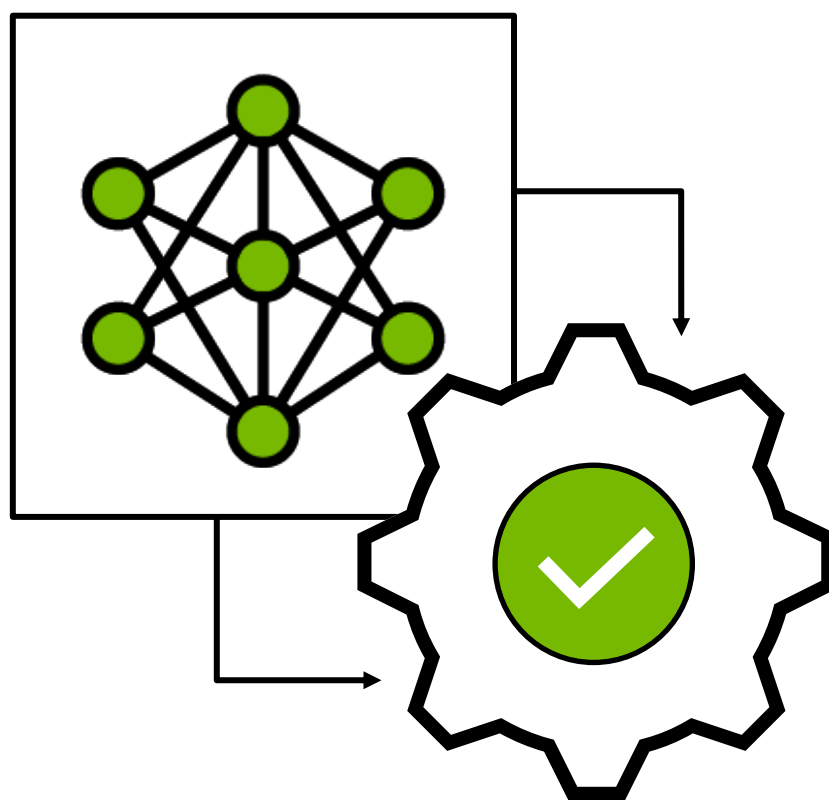1.5X per year

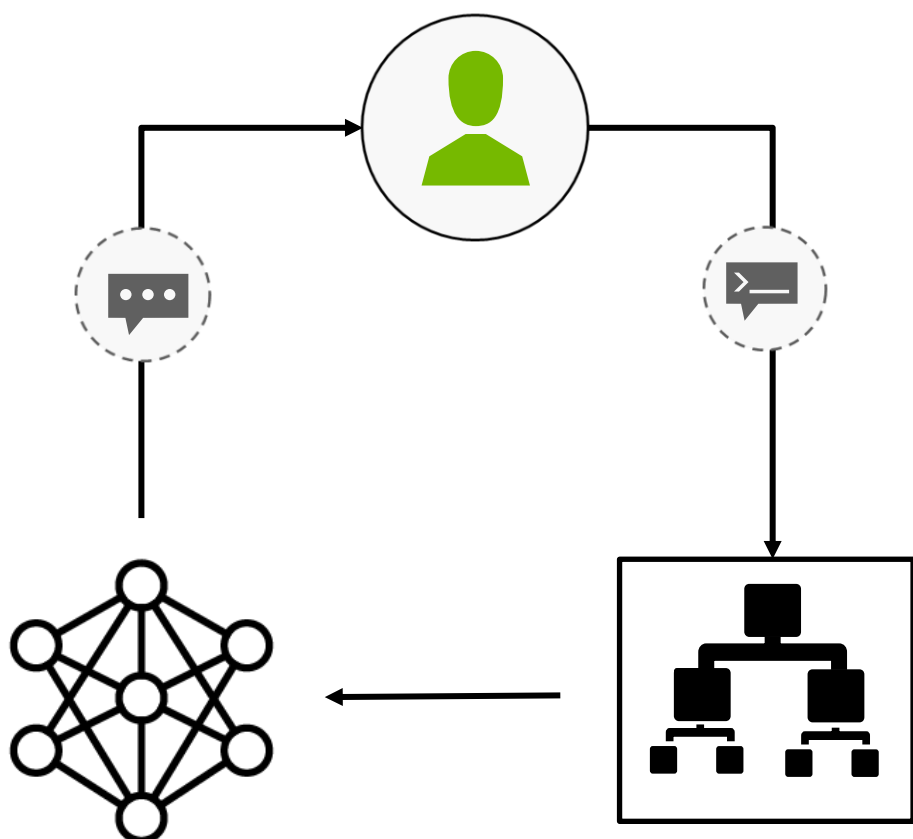Single-threaded perf

1980    1990    2000    2010    2020

# Data Processing for Different LLM Needs



Training Foundation Model

Fine-Tuning Foundation Model

Retrieval Augmented Generation (RAG)

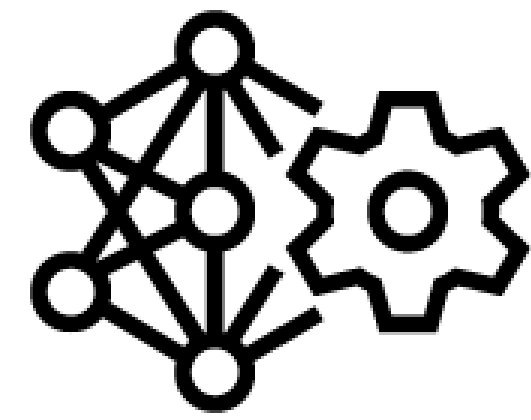| | Training Foundation Model | Fine-Tuning Foundation Model | Retrieval Augmented Generation (RAG) |
|---|---|---|---|
| Data Size | TB and PB | GBs | GBs |
| Compute Scale | Supercomputer | Single-node | Single GPU |
| Frequency | One-time | Iterative | Iterative & Continuous |

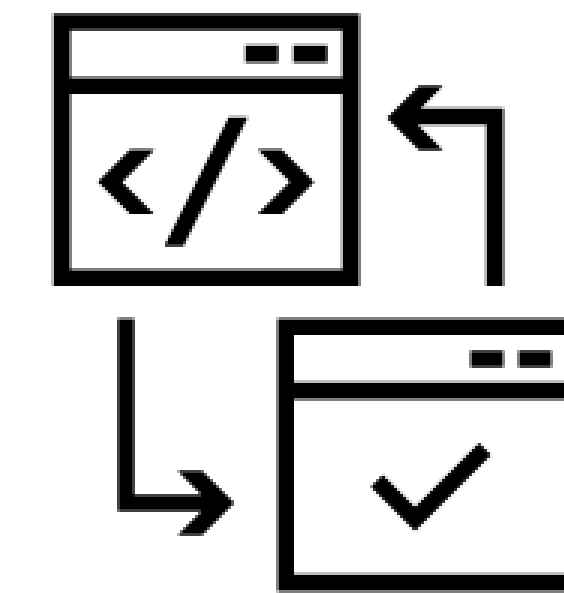# Challenges with Existing Solutions for Training Foundation Models

Inefficiencies lead to higher TCO and slower time to market

Longer Processing Time

Un-optimized Models

Un-optimized Pipelines

Knowledge & Expertise

# NeMo Curator - Overview

Scalable, configurable pipelines to curate text, image and video datasets lead to more accurate applications

| Higher Accuracy | Faster Processing | Scalability | Classifier Models | Deploy anywhere |
|---|---|---|---|---|

Improve accuracy with less data and less training compute

GPU acceleration with RAPIDS for **Dedupe (Exact, Fuzzy, Semantic)** and **Quality Classifier Models**
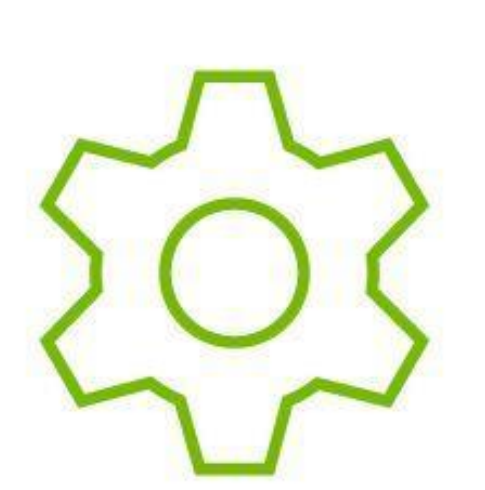
Up to **100+ PB data** by scaling across multiple nodes

State-of-the-Art quality classifier NIM microservice for safety, content, and diversity

Python APIs in a customizable and modular OSS library, runnable across CSP and On-prem

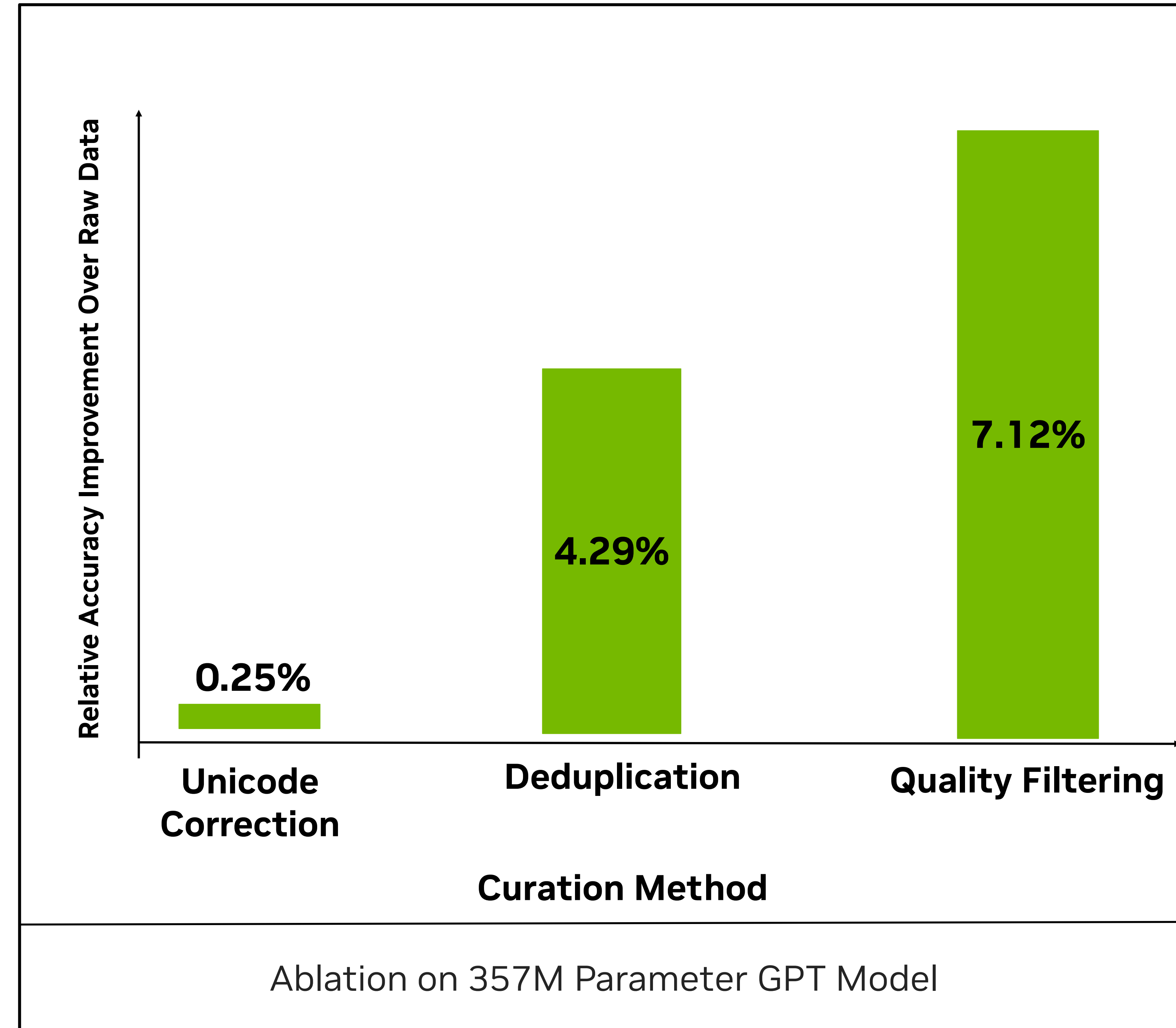| GitHub | NeMo framework container | PyPI | Microservice (coming soon) |
|---|---|---|---|

NVIDIA.

# Text Processing

# High Quality Data Processing Maximizes Model Performance

## Data Curation helps build SOTA Models

### LLM Accuracy Improvement on Curated Data



Relative Accuracy Improvement Over Raw Data

- Unicode Correction: 0.25%
- Deduplication: 4.29%
- Quality Filtering: 7.12%

Curation Method

Ablation on 357M Parameter GPT Model

# Why is Data Curation Important?

State-of-the-art data curation is essential for developing state-of-the-art models across all modalities

| Higher Accuracy | TCO Savings | Faster Training | Task Specialization |
|---|---|---|---|

Properly curated data leads to more improved accuracy across tasks
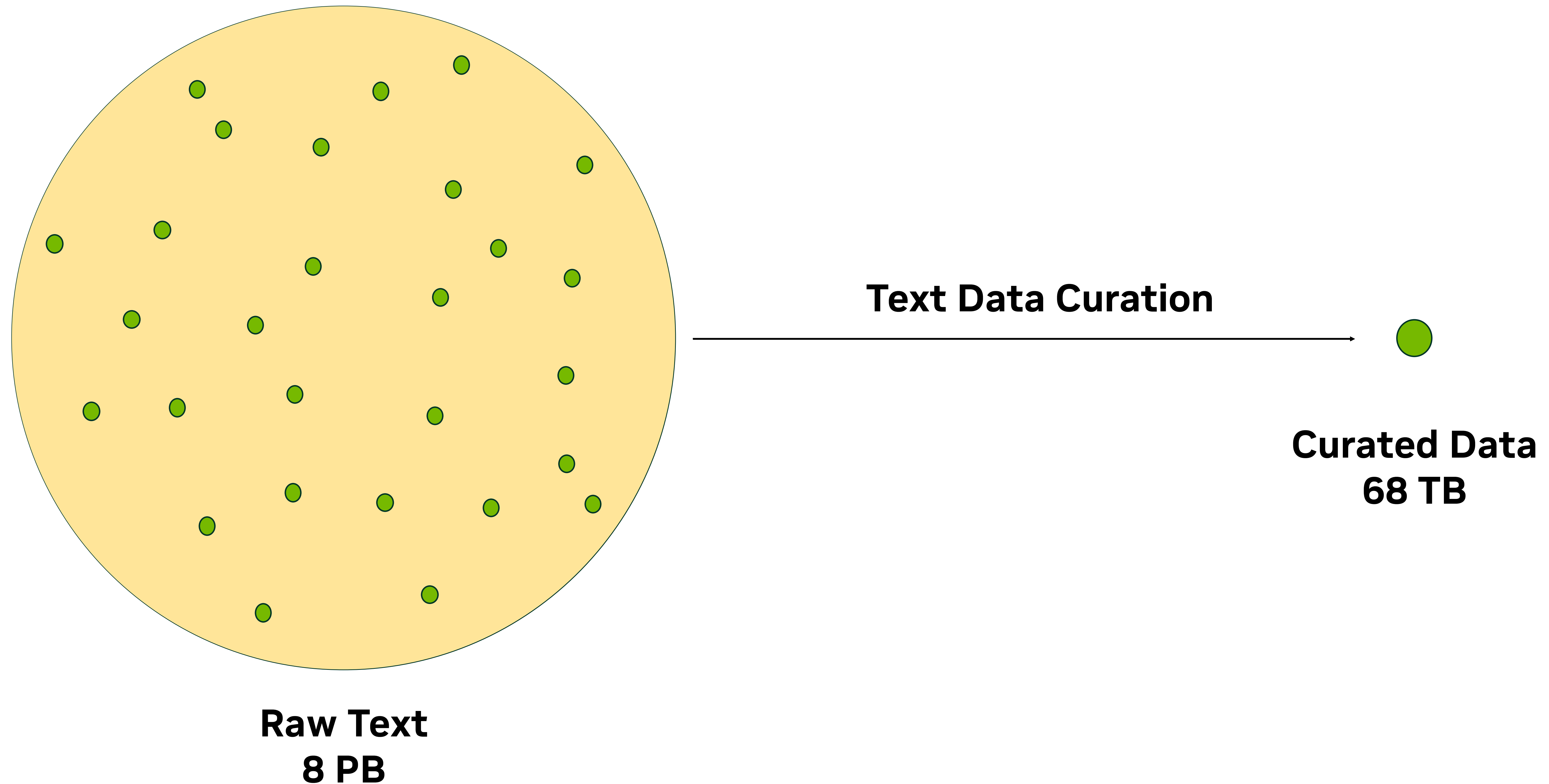
Decreases both training and inference costs significantly

Reduces compute requirements by orders of magnitude, enabling faster iterations
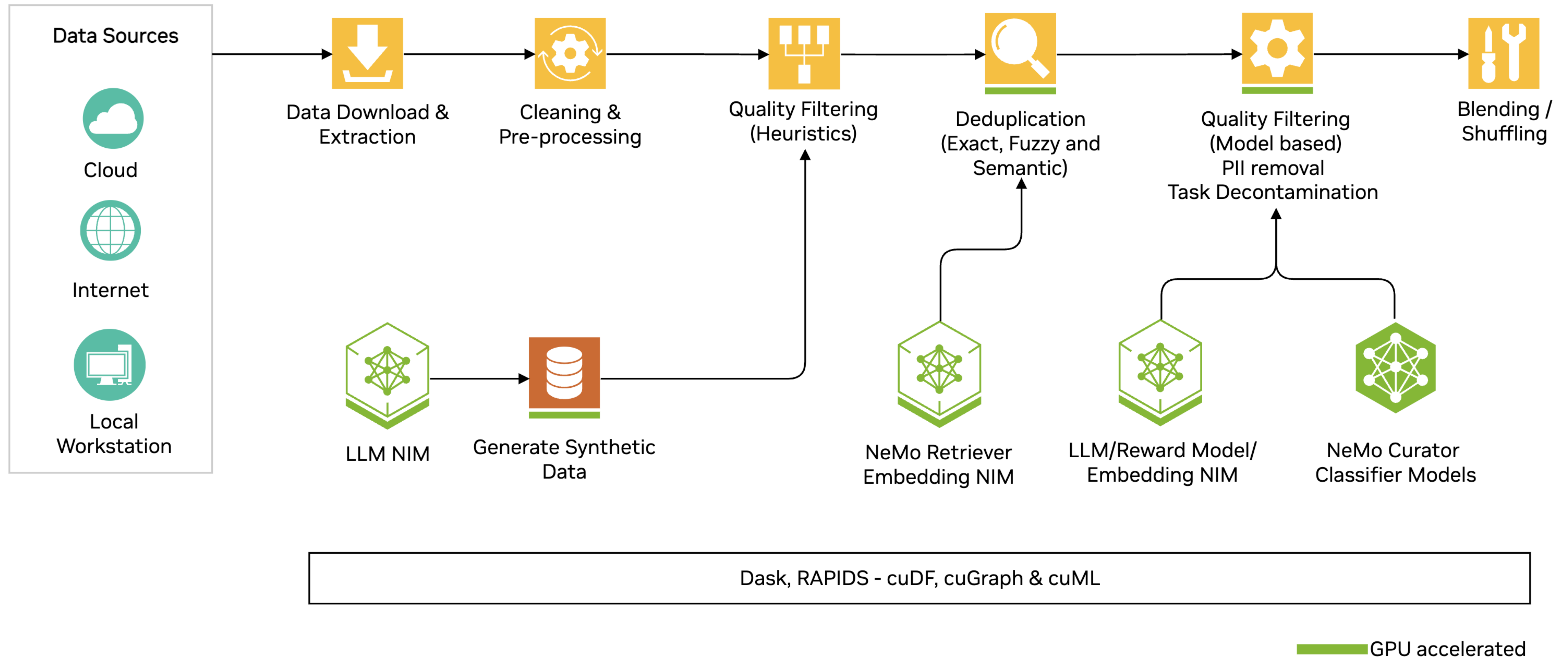
Enables optimization for specific tasks such as reasoning rather than general-purpose solutions

# Only 1% of Raw Text Data is Curated to Train Foundation Models



Text Data Curation

Raw Text
8 PB

Curated Data
68 TB

# Introducing: NeMo Curator for Text Processing
## Easily integrate different features into your existing pipelines with Python APIs



**Data Sources**

Cloud

Internet

Local Workstation

Data Download & Extraction → Cleaning & Pre-processing → Quality Filtering (Heuristics) → Deduplication (Exact, Fuzzy and Semantic) → Quality Filtering (Model based) PII removal Task Decontamination → Blending / Shuffling

LLM NIM → Generate Synthetic Data

NeMo Retriever Embedding NIM

LLM/Reward Model/ Embedding NIM

NeMo Curator Classifier Models

Dask, RAPIDS - cuDF, cuGraph & cuML

GPU accelerated

GitHub | NeMo framework container | PyPI

NVIDIA.

# NeMo Curator: Features to Train Foundation Models

## Achieve higher accuracy with a variety of GPU-accelerated features

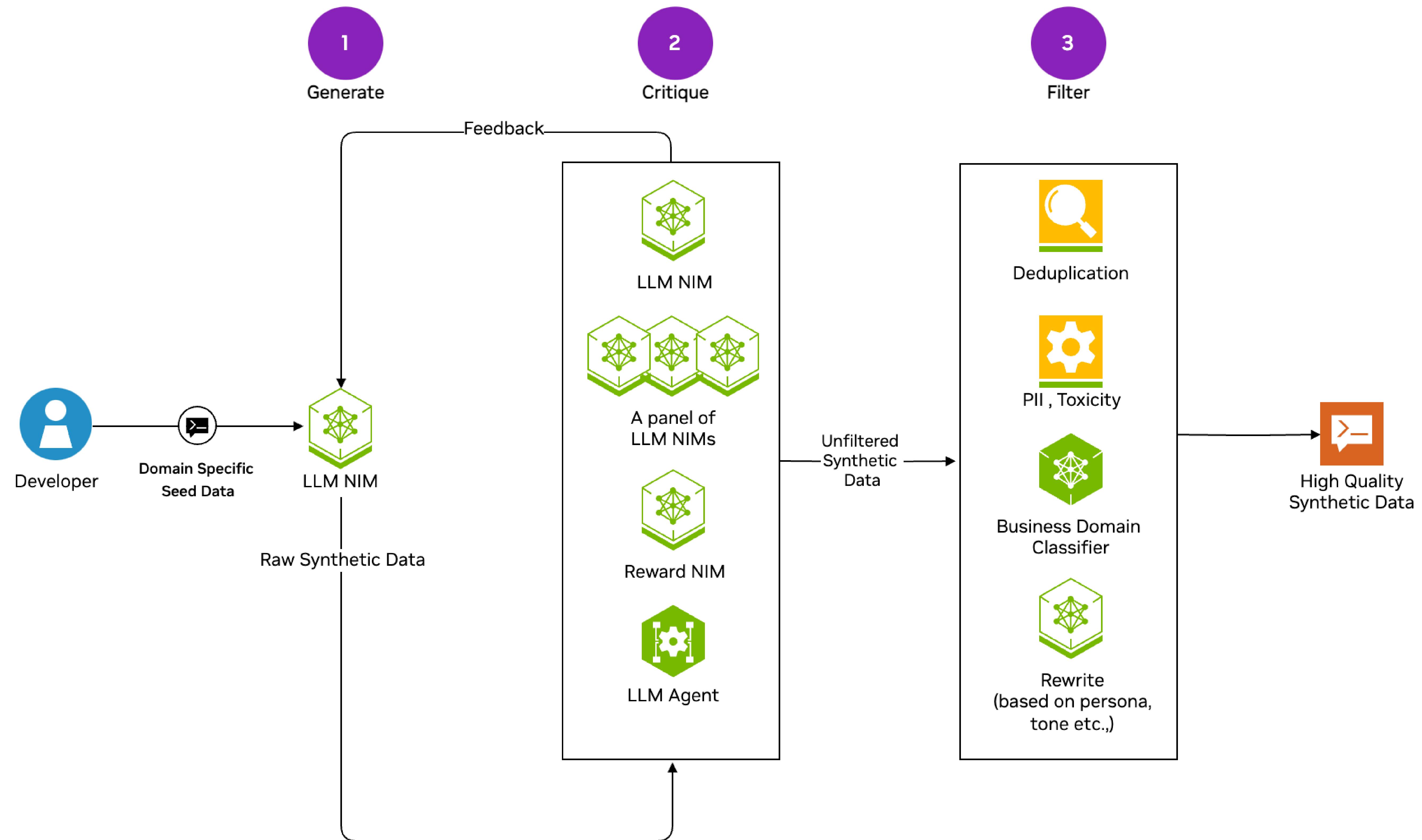| Synthetic Data Generation | Deduplication & Classification | GPU Acceleration with RAPIDS |
|---|---|---|

### Synthetic Data Generation

➤ **Pre-built pipelines** - for tasks like prompt generation, dialogue generation, and entity classification

➤ **Modular** - Easily integrate NeMo Curator's features into your existing pipelines

➤ **OpenAI API compatible** - Integrate custom Instruct and Reward models

### Deduplication & Classification

➤ **Lexical Deduplication** – Identical (Exact) or near identical (Fuzzy)

➤ **Semantic Deduplication** – focuses on the meaning rather than the exact text

➤ **Classifier Models** - State-of-the-art open models to either enrich or filter your data.

### GPU Acceleration with RAPIDS

➤ **cuDF** - for deduplication & classifer models

➤ **cuML** - for K-means clustering in semantic deduplication

➤ **cuGraph** – for fuzzy deduplication

# Synthetic Data Generation

Easily get started with pre-built pipelines or integrate various features into your existing workflows

# Synthetic Data Generation

Easily get started with pre-built pipelines or integrate various features into your existing workflows

| Pre-built Pipelines | Llama 3.1 Nemotron 70B/340B Reward Model | Tooling Support |
|---|---|---|
| • Prompt generation (open q/a, closed q/a, writing, math/coding )<br><br>• Synthetic two-turn prompt generation<br>  • Dialogue generation<br>  • Entity classification | • Benchmark-Topping Performance<br><br>• Single scalar scoring for human preferences<br><br>• Permissive license for commercial use | • Integrate into existing pipelines<br><br>• Bring your custom Instruct/Reward Model<br><br>• Supports different filtering techniques |

NVIDIA.

# Synthetic Data Generation: Example

## Use prompts to synthetically generate the QnAs and the reward score

```python
model = "nvdev/nvidia/llama-3.1-nemotron-70b-instruct"

question_lst = []

for i in range(len(ragas)):
    closed_qa_responses = client.query_model(
        model=model,
        messages = [
            {
                "role": "user",
                "content": f"Generate a single, concise question similar to '{ragas.question[i]}' based on the provided context.\n"
                            f"Do NOT include any explanations, rationale, or additional text in your response—just the question itself:\n"
                            f"{ragas.contexts[i]}"
            }
        ],
        temperature = 0.2,
        top_p = 0.7,
        max_tokens = 1024,
    )

    question = closed_qa_responses[0]
    question_lst.append(question)

ragas['sdg_nemotron_q'] = question_lst
```
✓ 8.7s

```python
ragas['sdg_nemotron_q']
```
✓ 0.0s

```
0    What is the significance of including a module...
1    How does the task execution stage in HuggingGP...
2    What are the primary technical hurdles encount...
3    How does Algorithm Distillation (AD) leverage ...
4    What are the core architectural elements and o...
5    How do consistent naming conventions, intra-fi...
6    How does API-Bank assess LLMs' decision-making...
7    How does a package.json fit into defining a No...
8    How does API-Bank assess LLMs' utilization of ...
9    How do finite context lengths and long-term pl...
Name: sdg_nemotron_q, dtype: object
```

```python
reward_lst = []

model = "nvdev/nvidia/llama-3.1-nemotron-70b-reward"

for i in range(len(ragas)):
    messages = [
        {
            "role": "user",
            "content": f"""
                Provide a concise and well-reasoned answer to the following question:
                "{ragas.sdg_nemotron_q[i]}"
                Use the context below to ensure accuracy and clarity:
                {ragas.contexts[i]}
                """
        },
        {
            "role": "assistant",
            "content": ragas.sdg_nemotron_a[i],
        },
    ]

    rewards = client.query_reward_model(messages=messages, model=model)
    rewards = float(rewards)

    reward_lst.append(rewards)

ragas['sdg_nemotron_reward'] = reward_lst
```
✓ 3.1s

```python
ragas['sdg_nemotron_reward']
```
✓ 0.0s

```
0    -13.93750
1     -4.09375
2     -8.68750
3     -8.06250
4     -5.15625
5    -15.75000
6     -8.18750
7     -9.31250
8     -7.56250
9    -11.06250
Name: sdg_nemotron_reward, dtype: float64
```

**NVIDIA**

# SDG Demo Video

# Synthetic Data Generation: Example
### Build an SDG pipeline to generate the data

## main.py

```
14
15  > import argparse~
38
39  SCRIPT_DIR_PATH = os.path.dirname(os.path.abspath(__file__))
40  DATA_DIR = os.path.join(SCRIPT_DIR_PATH, "data")
41  TEMP_DIR = os.path.join(SCRIPT_DIR_PATH, "_temp")
42  CONFIG_DIR = os.path.join(SCRIPT_DIR_PATH, "config")
43  DATASET_URL = "https://huggingface.co/datasets/ymoslem/Law-StackExchange/resolve/main/law-stackexchange-questions-answers.json"
44
45
46  > def pre_imports():~
48
49
50  > def random_split_rows(rows: List[Any], train_ratio: float, val_ratio: float, seed=42):~
72
73
74  > def download_and_convert_to_jsonl() -> str:~
114
115
116 > def semantic_dedupe(dataset):~
144
145
146 > def run_curation_pipeline(~
237
238
239 > def run_pipeline(args, jsonl_fp):~
364
365
366 > def main():~
432
433
434 > if __name__ == "__main__":~
```

## syn_gen.py

```
14
15  > import asyncio~
26
27  > PROMPT_GENERATE_QUESTIONS_FROM_ANSWER = """TEXT:~
36
37  > PROMPT_PARAPHRASE_TEXT = """TEXT:~
46
47
48   class SyntheticGenerator:
49
50  >     def __init__(~
81
82  >     def run(~
106
107 >     def _split_sdg_responses(self, sdg_response: str) -> List[str]:~
120
121 >     def _write_all_to_file(self, gen_entries, out_fp: str):~
182
183 >     async def _prompt_model(~
252
253 >     async def _synthesize_from_source(~
```

## multiple rounds of SDG

```
∨ data
  ∨ curated
    ∨ final
      {} law-qa-test.jsonl
      {} law-qa-train.jsonl
      {} law-qa-val.jsonl
    ∨ round-1
      {} law-qa-train-synth-round-1.jsonl
      {} law-qa-train.jsonl
    > round-2
    > round-3
    > round-4
    > round-5
  ∨ raw
    ∨ downloads
      {} law-stackexchange-questions-answers.json
    ∨ splits
      {} law-qa-test.jsonl
      {} law-qa-train.jsonl
      {} law-qa-val.jsonl
```

## one example on SDG curated result

```
{
    "id": "appliance-qa-B00074TB9U-synth-0",
    "asin": "B00074TB9U-synth-0",
    "question": "How effective is a ductless installation in eliminating cooking smells?",
    "answer": "The vent has two speeds; the lower one is suitable for regular air removal, while the higher speed is necessary for heavy-duty cooking. The higher speed functions well.",
    "questionType": "yes/no",
    "score": -1,
    "helpfulness": 0.6015625,
    "correctness": 0.388671875,
    "coherence": 3.078125,
    "complexity": 0.8671875,
    "verbosity": 0.60546875
}
```

NVIDIA.

# Deduplication and Filtering

## Scale across multi-node, multi-GPU setups, eliminating the need for iterative CPU processing

- Supports lexical and semantic deduplication for document processing

- Scales on multi node, multi-GPU by leveraging RAPIDS
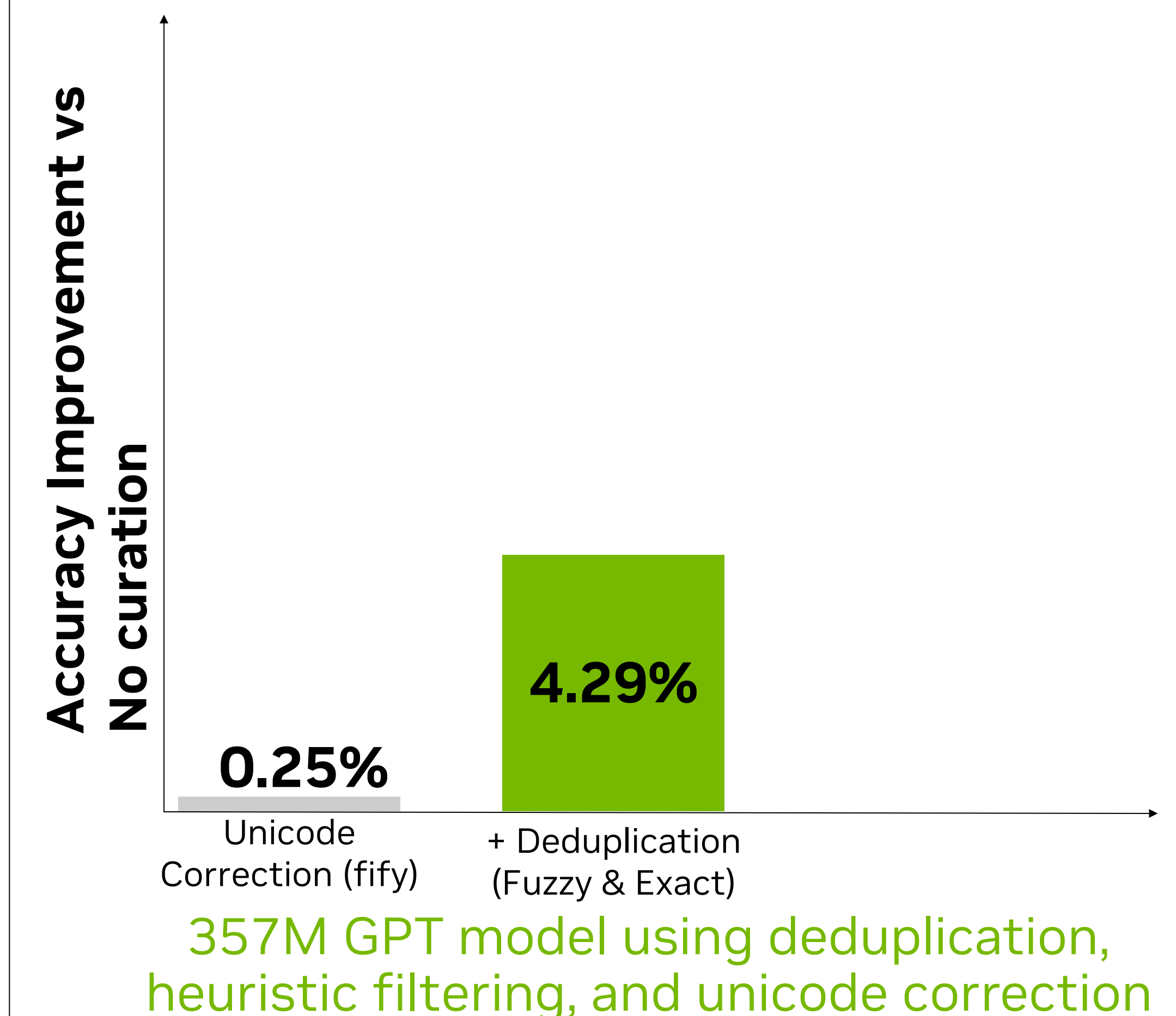
- Scale to 100+ TB of data

| Lexical Deduplication | Semantic Deduplication | Acceleration |
|---|---|---|

**Lexical Deduplication**

- Utilizes text similarity to discover duplication

- **Exact deduplication**
  - Hashing based matching for each document

- **Fuzzy deduplication**
  - Minhash, Bucketization and Clustering

**Semantic Deduplication**

- Utilizes meaning of text to discover duplication

- Leverages embedding models to identify semantic documents

- Modularity to use custom embedding model

**Acceleration**

Accuracy Improvement vs No curation

**0.25%**
Unicode Correction (fify)

**4.29%**
+ Deduplication (Fuzzy & Exact)

357M GPT model using deduplication, heuristic filtering, and unicode correction

NVIDIA.

# NeMo Curator - Deduplication

Scale across multi-node, multi-GPU setups, eliminating the need for iterative CPU processing

## Value Proposition

➢ Accelerate on multi-node multi-GPU setups by leverage RAPIDS

➢ Scale up to 100+ TBs of data

➢ Get support for both

   ➢ Lexical Deduplication: Uses text similarity, both Exact and Fuzzy

   ➢ Semantic Deduplication: Uses text meaning

## Technical Details

### Lexical Deduplication

➢ Exact deduplication

   ➢ Hashing for each document

➢ Fuzzy deduplication

   ➢ Minhash, Bucketization and Clustering

### Semantic Deduplication

➢ Leverage embeddings to identify semantic documents

➢ Use SOTA NeMo Retriever Text Embedding NIM or customize with your embedding model

NVIDIA

# NeMo Curator - Classifier Models

## Create high-quality data blends with RAPIDS accelerated inference

## Value Proposition

➢ Accelerated inference with RAPIDS powered distributed data classification module and intelligent batching

➢ Seamless scalability for classifying TBs of data

➢ Faster processing with parallelization across multiple GPUs

➢ Lower memory and compute footprint with state-of-the-art open classifier models (Apache 2.0 license )

➢ 8 classifier models released

## Classifier Models

### Domain Classifier

➢ Supports 26 domain classes

    ➢ Top 10 classes : Finance, Health, Business and Industrial, Science, Law and Government, Internet and Telecom, Jobs and Education, News, Computers and Electronics, Shopping;

➢ Trained on 1 million Common Crawl samples & 500k Wikipedia articles

➢ Also available for multiple languages

### Quality Classifier

➢ Classify quality of the document to 'High', 'Medium' or 'Low'

➢ Training data annotated by humans on quality factors such as: content accuracy, clarity, coherence, grammar, depth of information and overall usefulness of the document

# NeMo Curator - Classifier Models

## Create high-quality data blends with RAPIDS accelerated inference

### Prompt Task and Complexity Classifier

➢ Classify tasks across 11 categories

  ➢ OpenQA, Closed QA, Summarization, Text Generation, Code Generation, Chatbot, Classification, Rewrite, Brainstorming, Extraction, Other

➢ Compute complexity on 6 dimensions

  ➢ Creativity, Domain Knowledge, Reasoning, Constraints, Contextual Knowledge, # of Few shots

### FineWeb Nemotron-4 Edu Classifier

➢ Determine the educational value of a piece of text (score 0-5 from low to high).

➢ Trained using annotations from Nemotron-4-340B-Instruct
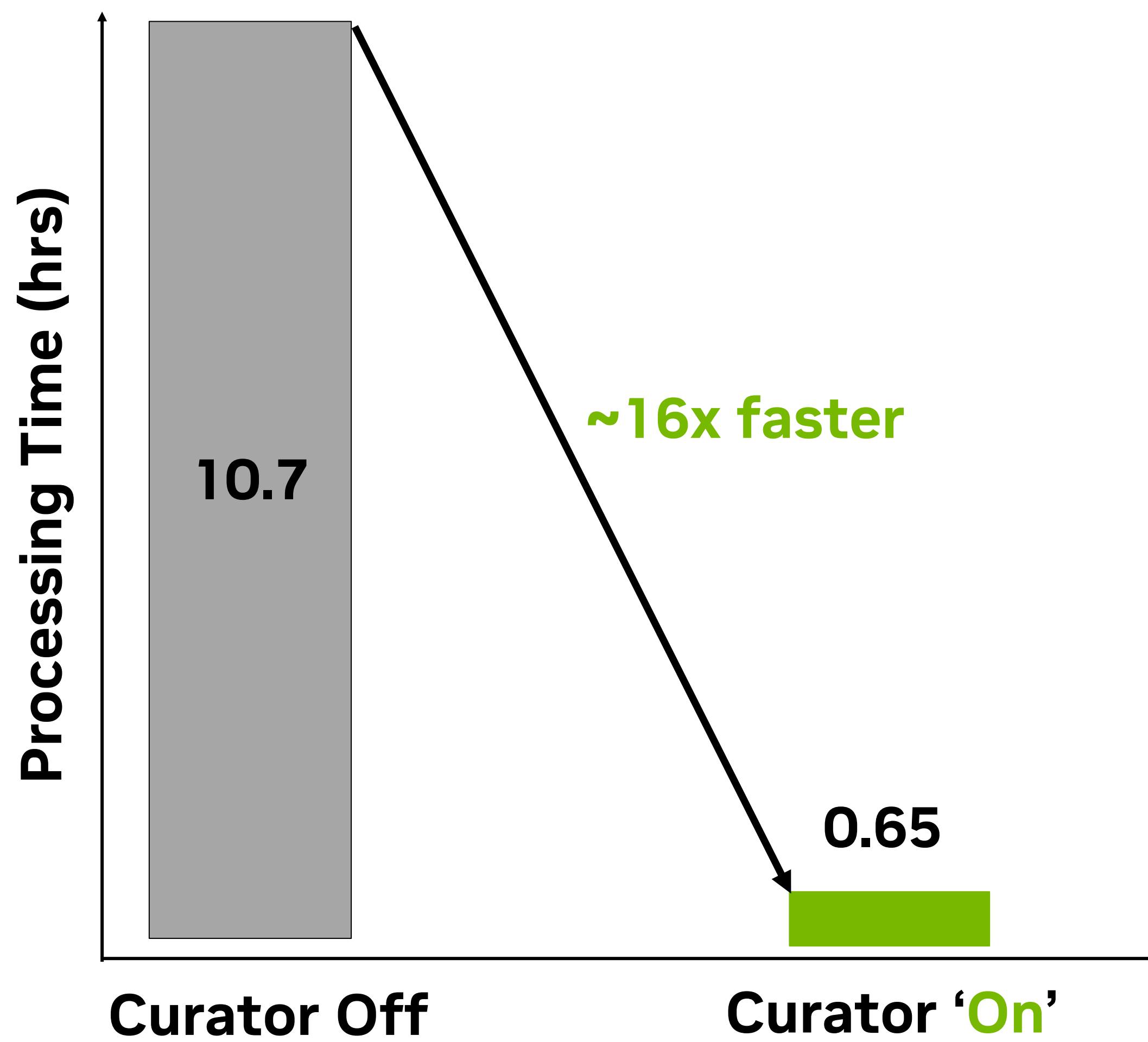
### Content Type Classifier

➢ Categorize documents into 11 content types

  ➢ Explanatory articles, News, Blogs, Boilerplate content, Analytical exposition, Online Comments, Reviews, Books & literature, Conversational, and Personal Websites.

➢ Useful for content management systems, digital publishers, recommendation systems

### Instruction Data Guard

➢ Identify malicious prompts used for fine-tuning

➢ Optimal use for instruction:response datasets
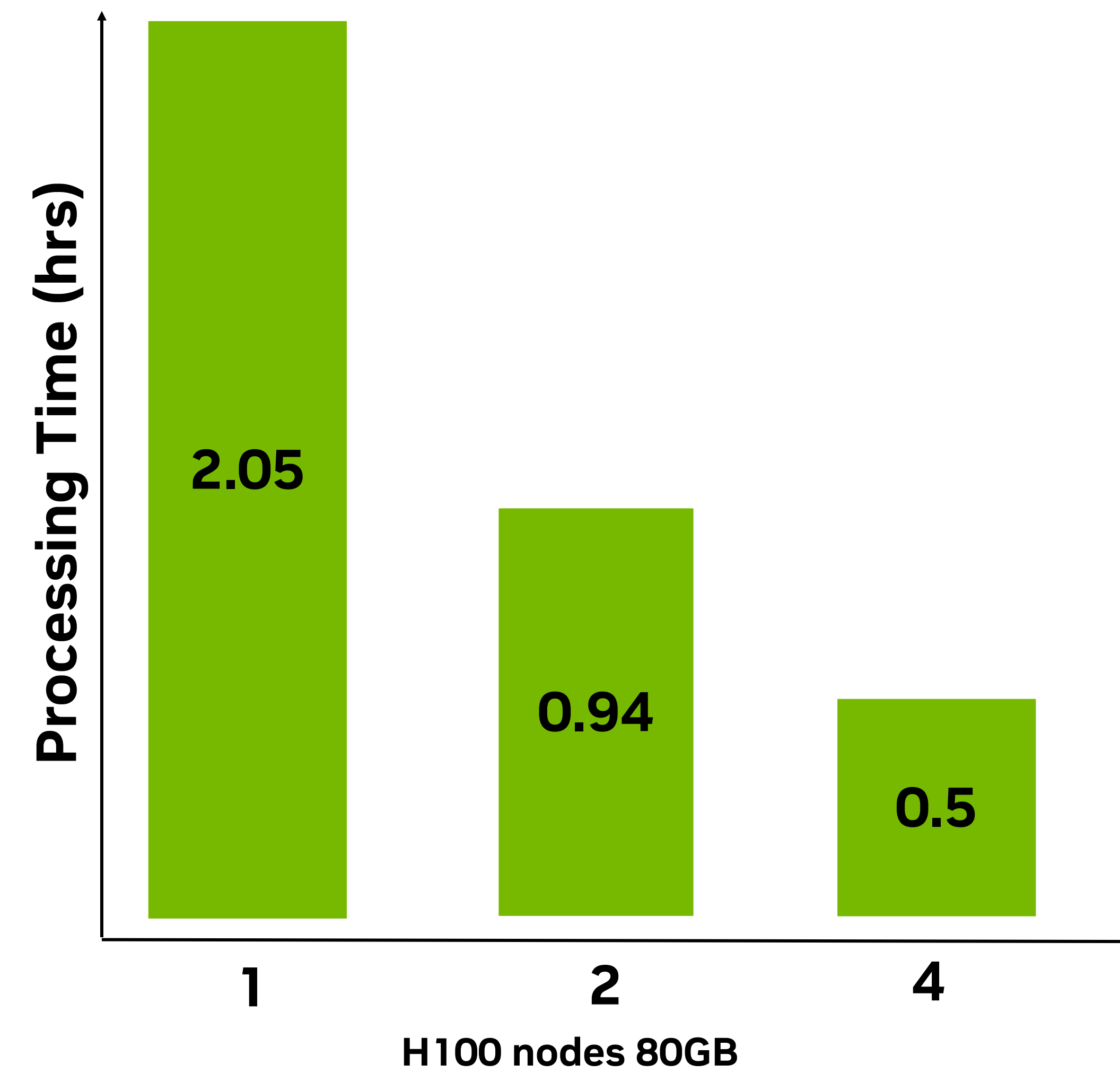
# Performance – Fuzzy Deduplication

**~16x Faster Processing**



Processing time for fuzzy deduplication of RedPajama-v2 subset (8TB/1.78T tokens)

'On': Data processed with NeMo Curator on 3 H100 nodes

'Off' : Data processed with a leading alternative library on CPUs

Processing time for fuzzy deduplication of RedPajama-v2 subset (8TB/1.78T tokens)

Scaling on 1, 2, 3, 4 H100 nodes 80GB

**NVIDIA.**

# Deep Dive: Fuzzy Deduplication

Identifying similar documents

Exact Duplicate: "We the People of the United States, in Order to form a more perfect Union, establish Justice..."

Document: "We the People of the United States, in Order to form a more perfect Union, establish Justice..."
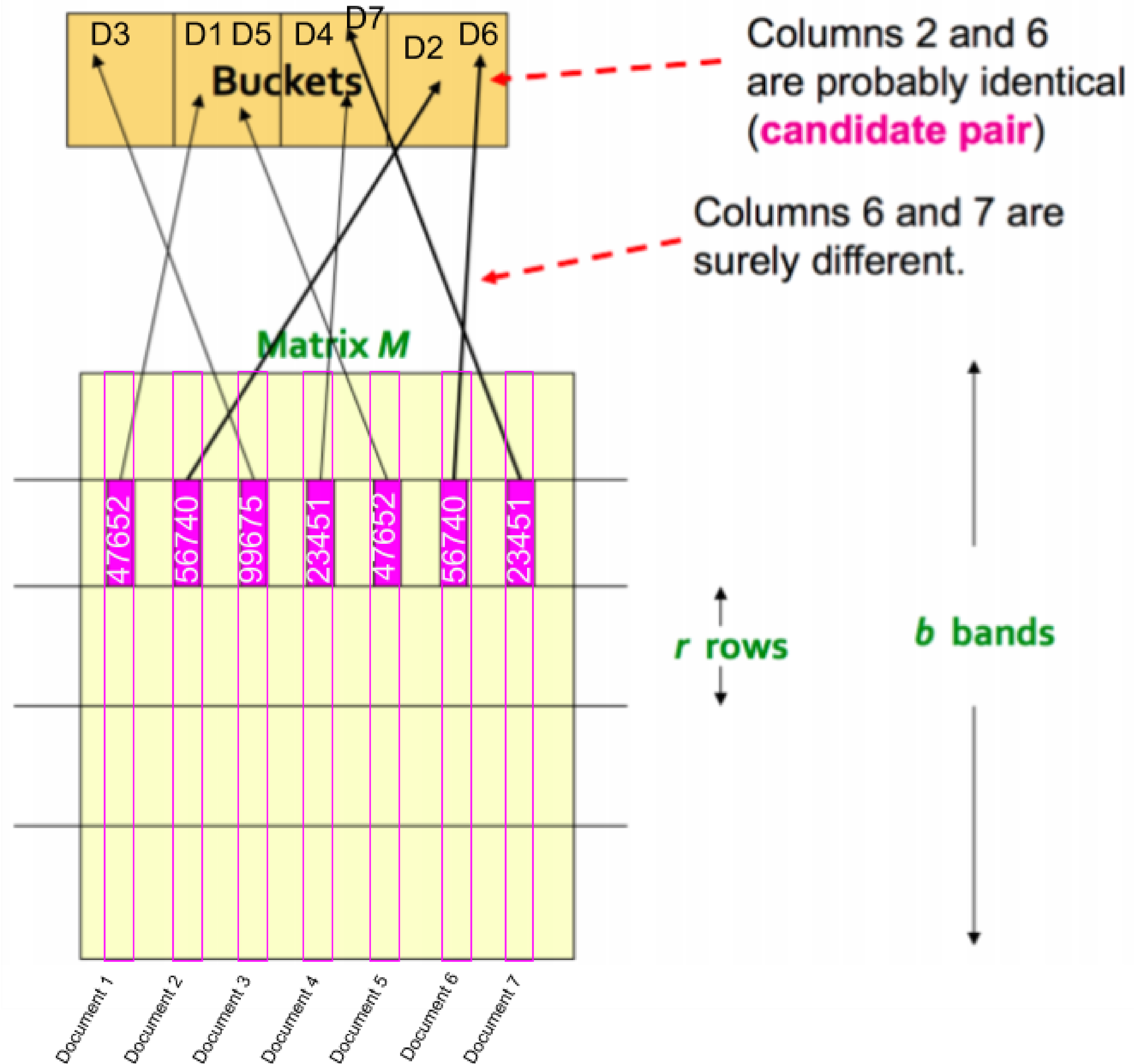
Fuzzy Duplicate: "Here is the US Constitution. We the People of the United States, in Order to form a more perfect Union, ensure Justice..."

Fuzzy Duplicate: "We the People of the United States, in Order to ensure Justice and domestic tranquility..."

NVIDIA.

# Min-hashing

| Doc -1 | | | | |
|---|---|---|---|---|
| K-shingle | {The quick brown} | {fox jumps over} | {the lazy dog} | |
| hashed-shingle | 345L | 3455L | 934L | |
| Hash_1 | 23 | 49 | 50 | 23 |
| Hash_2 | 56 | 24 | 39 | 24 |
| Hash_3 | 38 | 56 | 84 | 38 |
| Hash_4 | 48 | 29 | 93 | 29 |
| Hash_5 | 67 | 75 | 59 | 59 |

**Signature-1**

| Doc -2 | | | | |
|---|---|---|---|---|
| K-shingle | ... | ... | ... | |
| hashed-shingle | ... | ... | ... | |
| Hash_1 | ... | ... | ... | 34 |
| Hash_2 | ... | ... | ... | 56 |
| Hash_3 | ... | ... | ... | 78 |
| Hash_4 | ... | ... | ... | 23 |
| Hash_5 | ... | ... | ... | 14 |

**Signature-2**

| Doc -3 | | | | |
|---|---|---|---|---|
| ... | | | | |

Src: https://mrhasankthse.github.io/riz/2020/03/19/Minhash-and-LSH.html

# Signatures over hash buckets
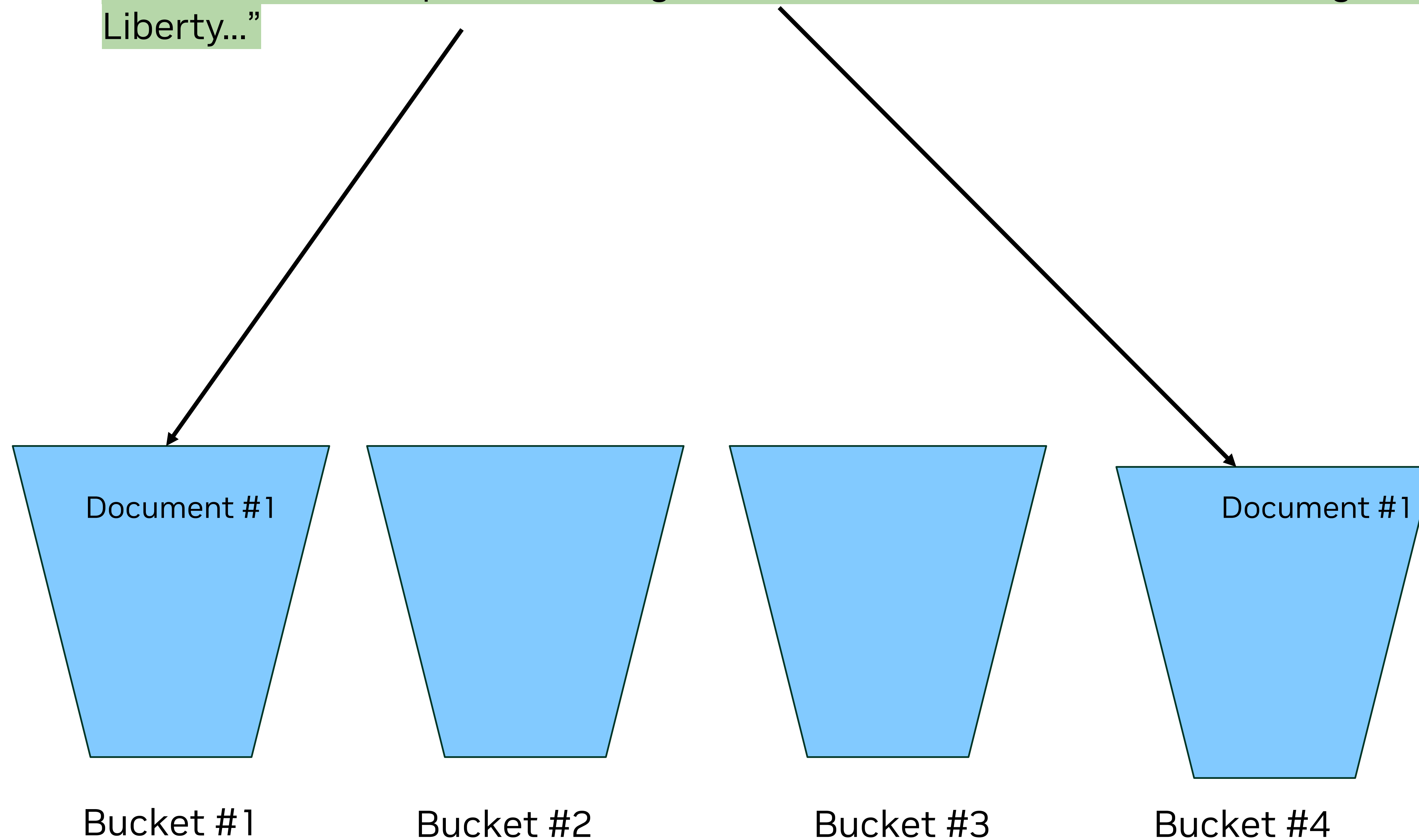


Src: https://mrhasankthse.github.io/riz/2020/03/19/Minhash-and-LSH.html

# Deep Dive: Fuzzy Deduplication
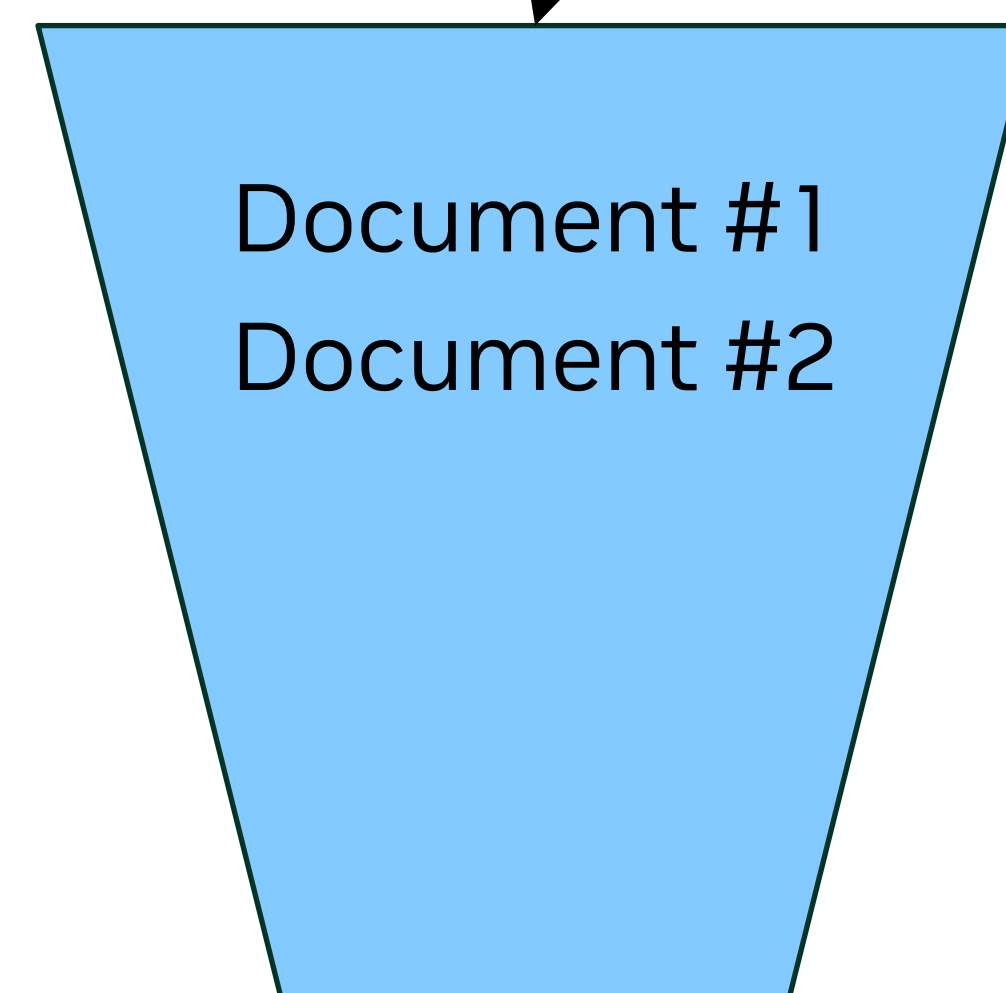## Min-Hashing and Bucketization

Document 1: "We the People of the United States, in Order to form a more perfect Union, establish Justice, ensure domestic Tranquility, provide for the common defense, promote the general Welfare, and secure the Blessings of Liberty…"

Document #1

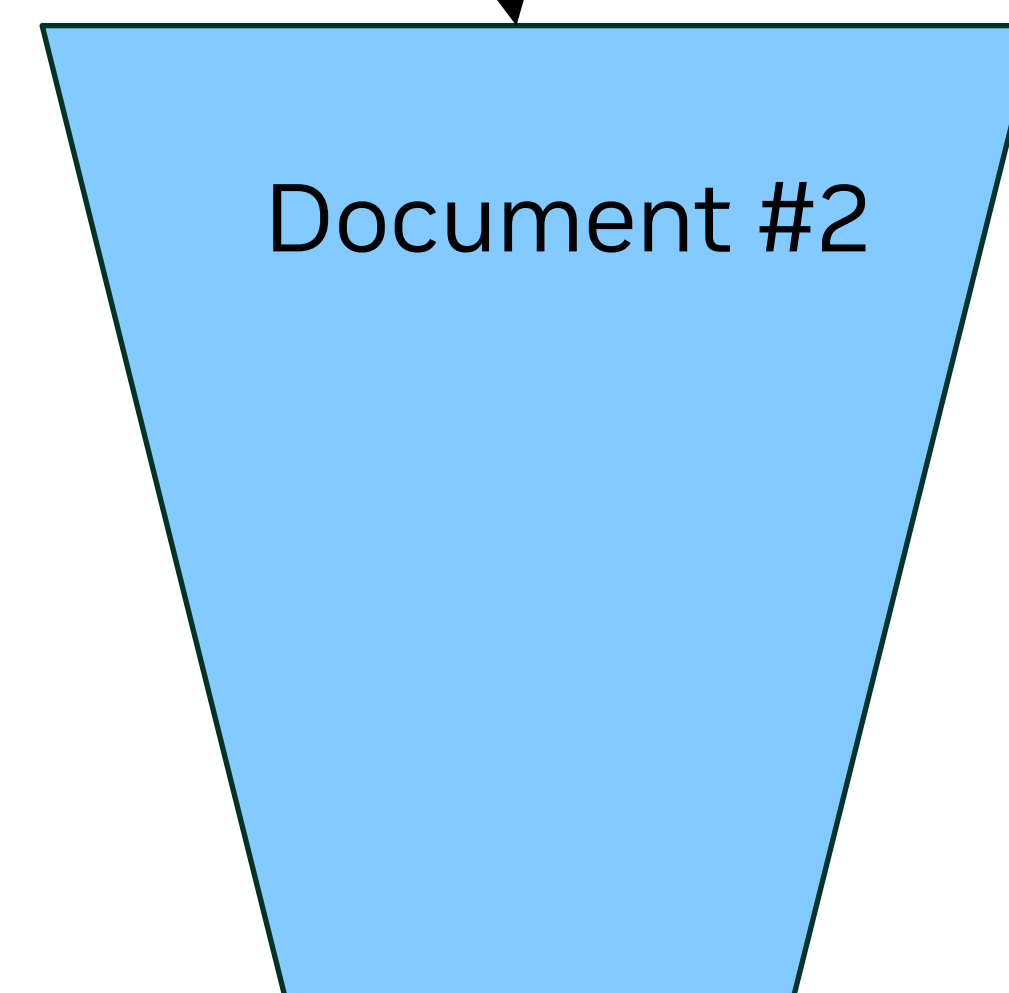Document #1

Bucket #1

Bucket #2

Bucket #3

Bucket #4

# Deep Dive: Fuzzy Deduplication
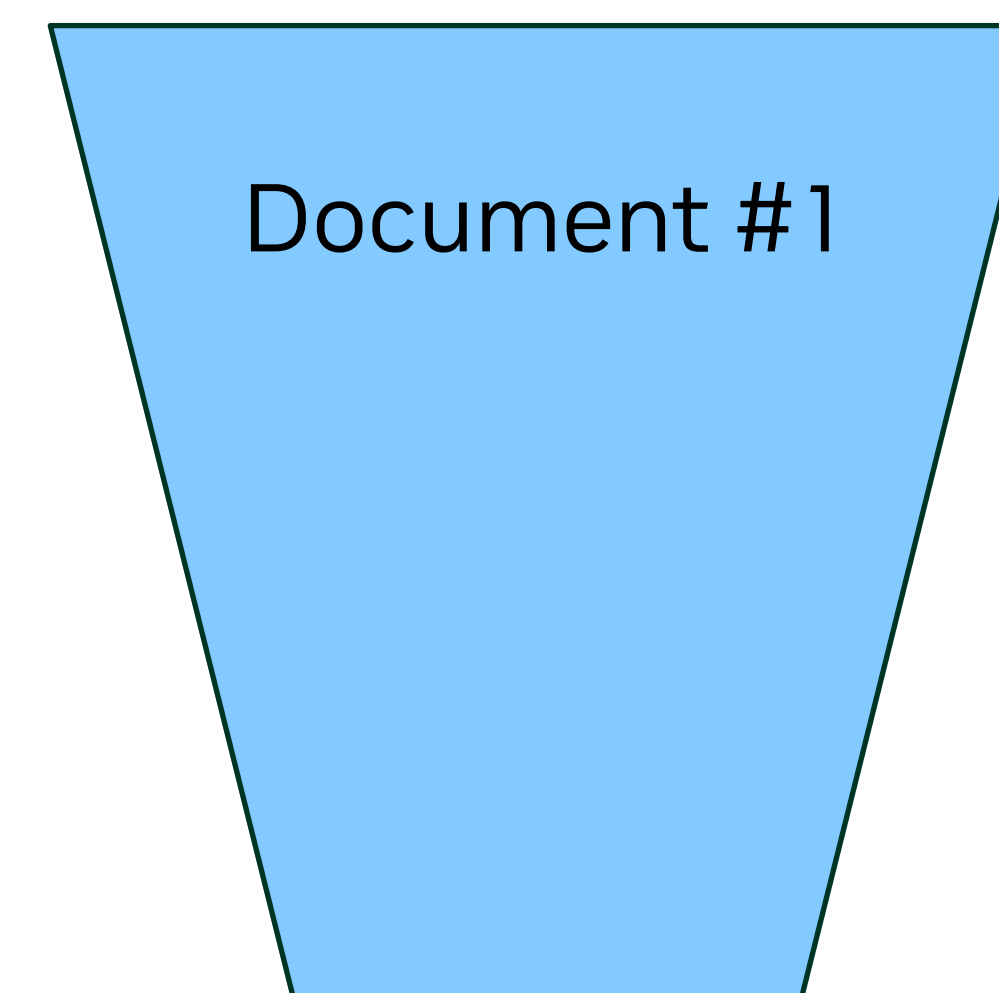Min-Hashing and Bucketization

Document 2: "Here is the US Constitution. We the People of the United States, in Order to form a more perfect Union, establish Justice, insure foreign Tranquility, provide for the common defense, promote the general Welfare, and secure the Blessings of Liberty..."
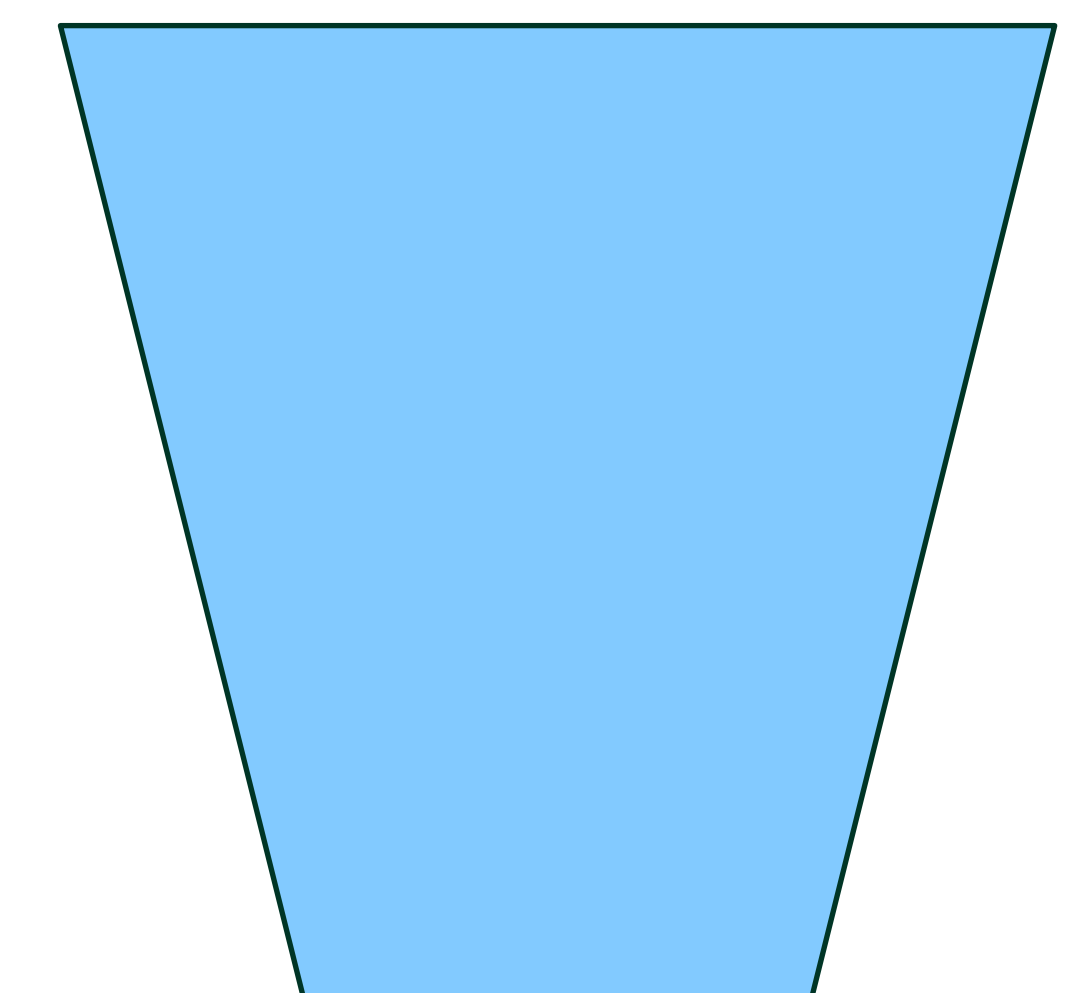
Document #1
Document #2

Document #2

Document #1

Bucket #1
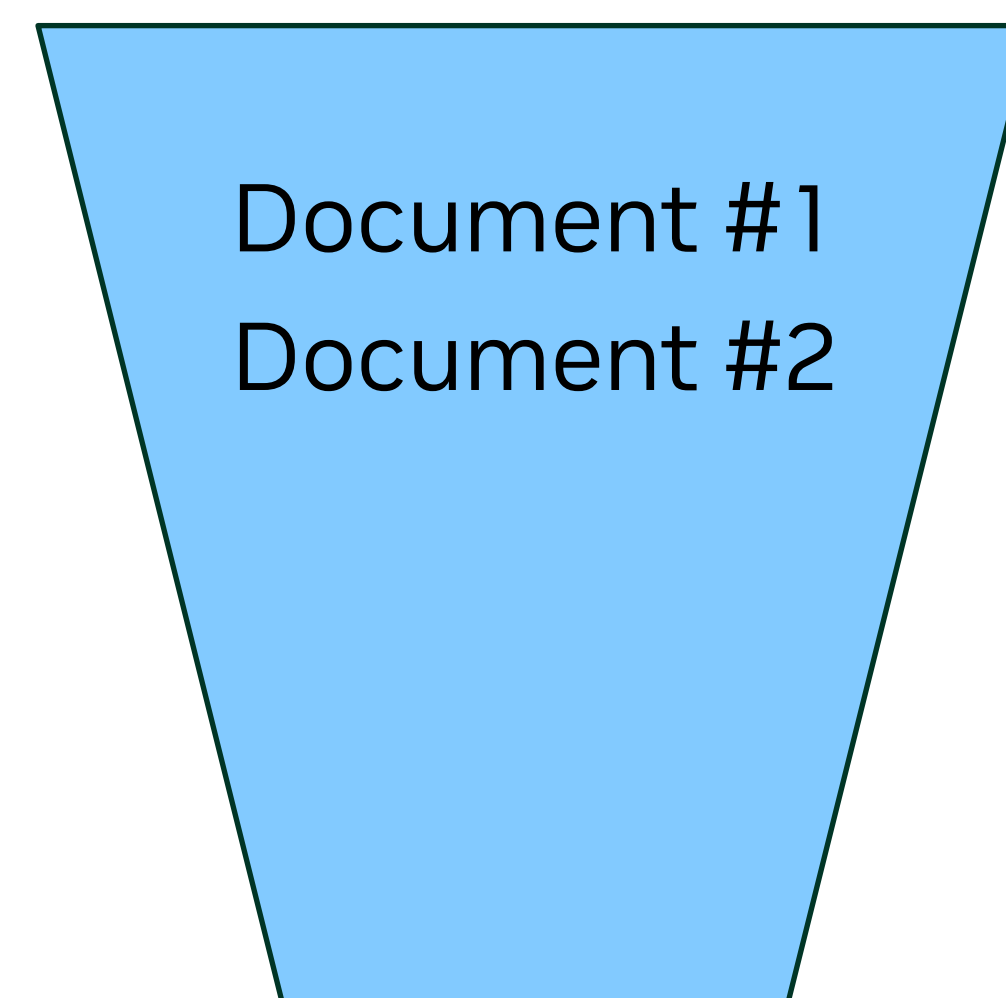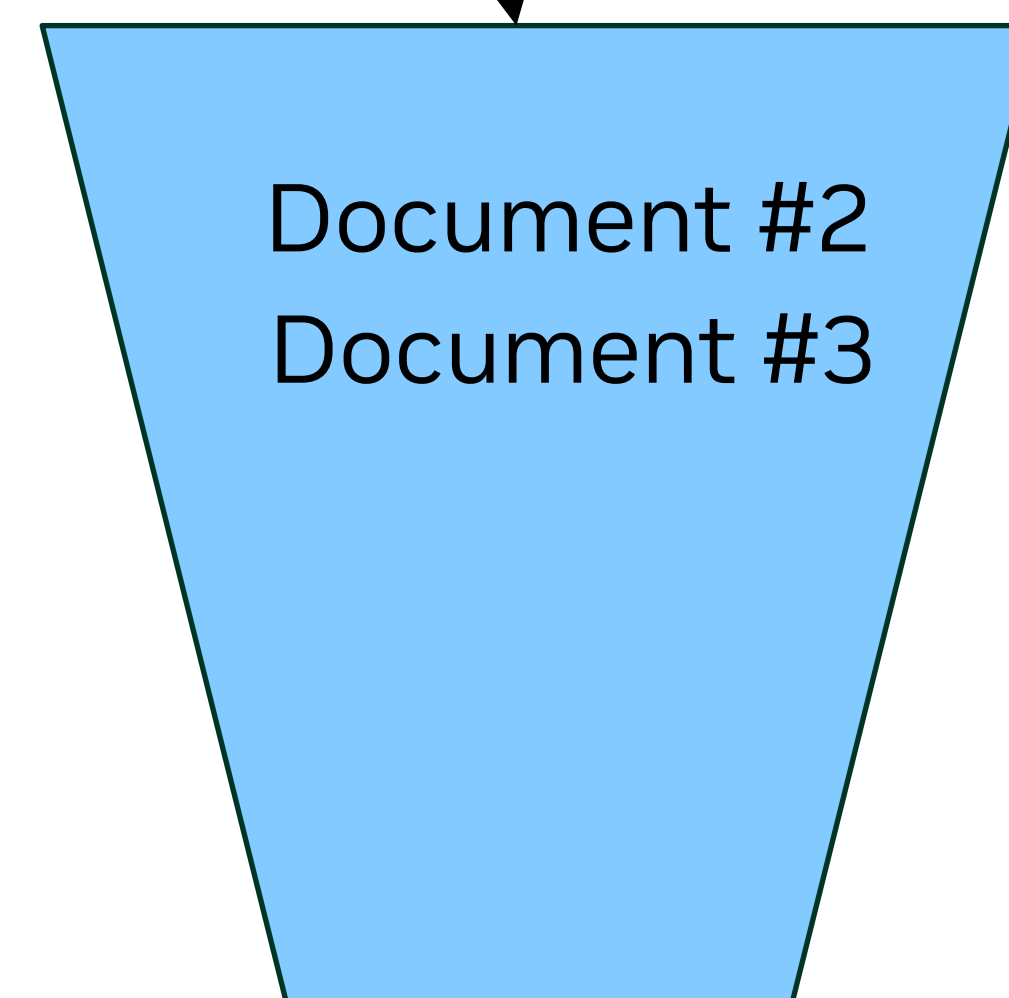
Bucket #2

Bucket #3

Bucket #4

# Deep Dive: Fuzzy Deduplication

## Min-Hashing and Bucketization

Document 3: "We the People of Canada, in Order to form a more perfect Union, establish Justice, insure domestic Tranquility, provide for the common defense, promote the general Welfare, and secure the Blessings of Liberty…"

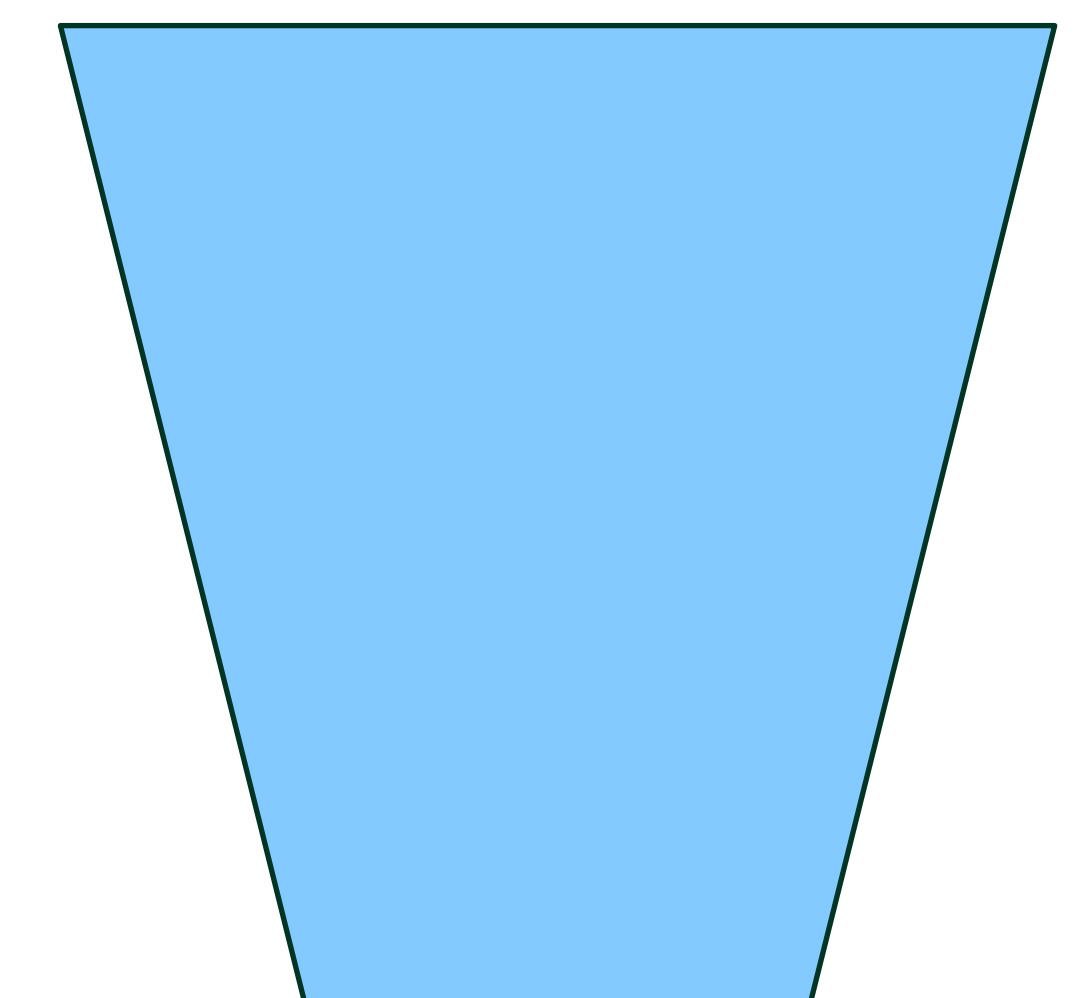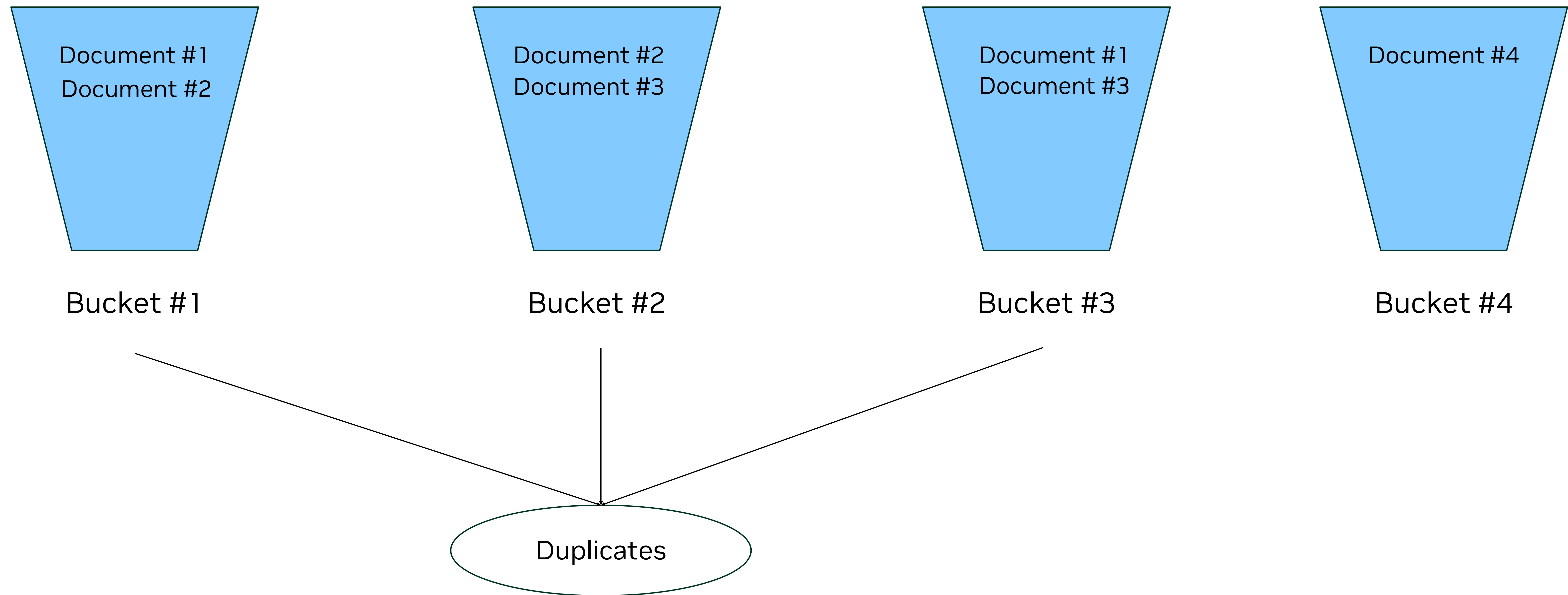| Bucket #1 | Bucket #2 | Bucket #3 | Bucket #4 |
|---|---|---|---|
| Document #1<br>Document #2 | Document #2<br>Document #3 | Document #1<br>Document #3 | |

# Deep Dive: Fuzzy Deduplication
## Min-Hashing and Bucketization

# Semantic Deduplication: Example

Remove redundant data by identifying and eliminating semantically similar data points

sem_dedup_config.yaml

```yaml
# Configuration file for semantic dedup
cache_dir: "semdedup_cache"
num_files: -1

# Embeddings configuration
embeddings_save_loc: "embeddings"
embedding_model_name_or_path: "sentence-transformers/all-MiniLM-L6-v2"
embedding_batch_size: 128

# Clustering configuration
clustering_save_loc: "clustering_results"
n_clusters: 1000
seed: 1234
max_iter: 100
kmeans_with_cos_dist: false

# Semdedup configuration
which_to_keep: "hard"
largest_cluster_size_to_process: 100000
sim_metric: "cosine"

# Extract dedup configuration
eps_thresholds:
  - 0.01
  - 0.001

# Which threshold to use for extracting deduped data
eps_to_extract: 0.01
```

NVIDIA

# Classification and PII
## Create high-quality data blends with RAPIDS accelerated inference

- Accelerated inference through distributed classification model and intelligent batching

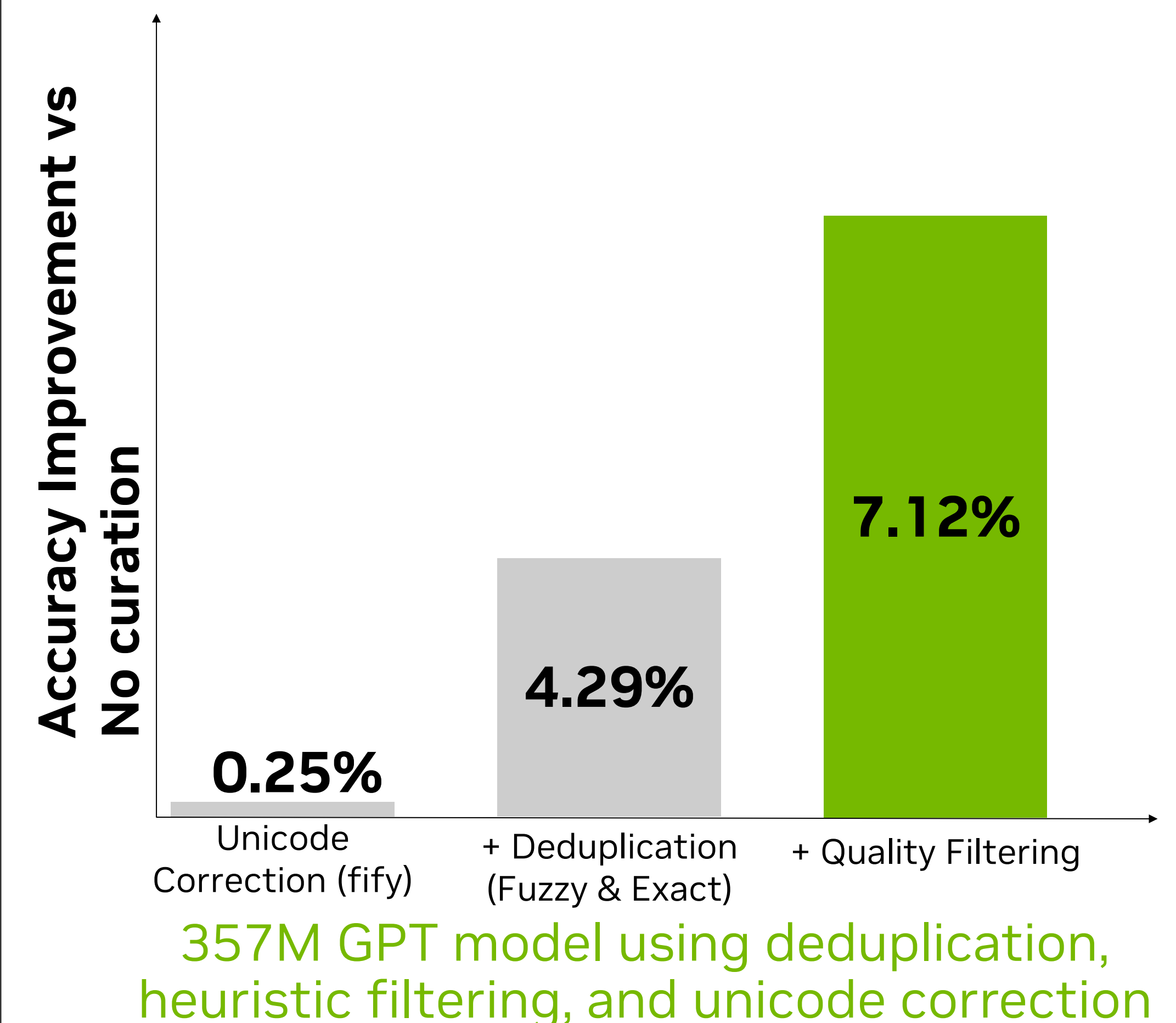- Redact/remove Personally Identifiable information (PII) using SOTA model

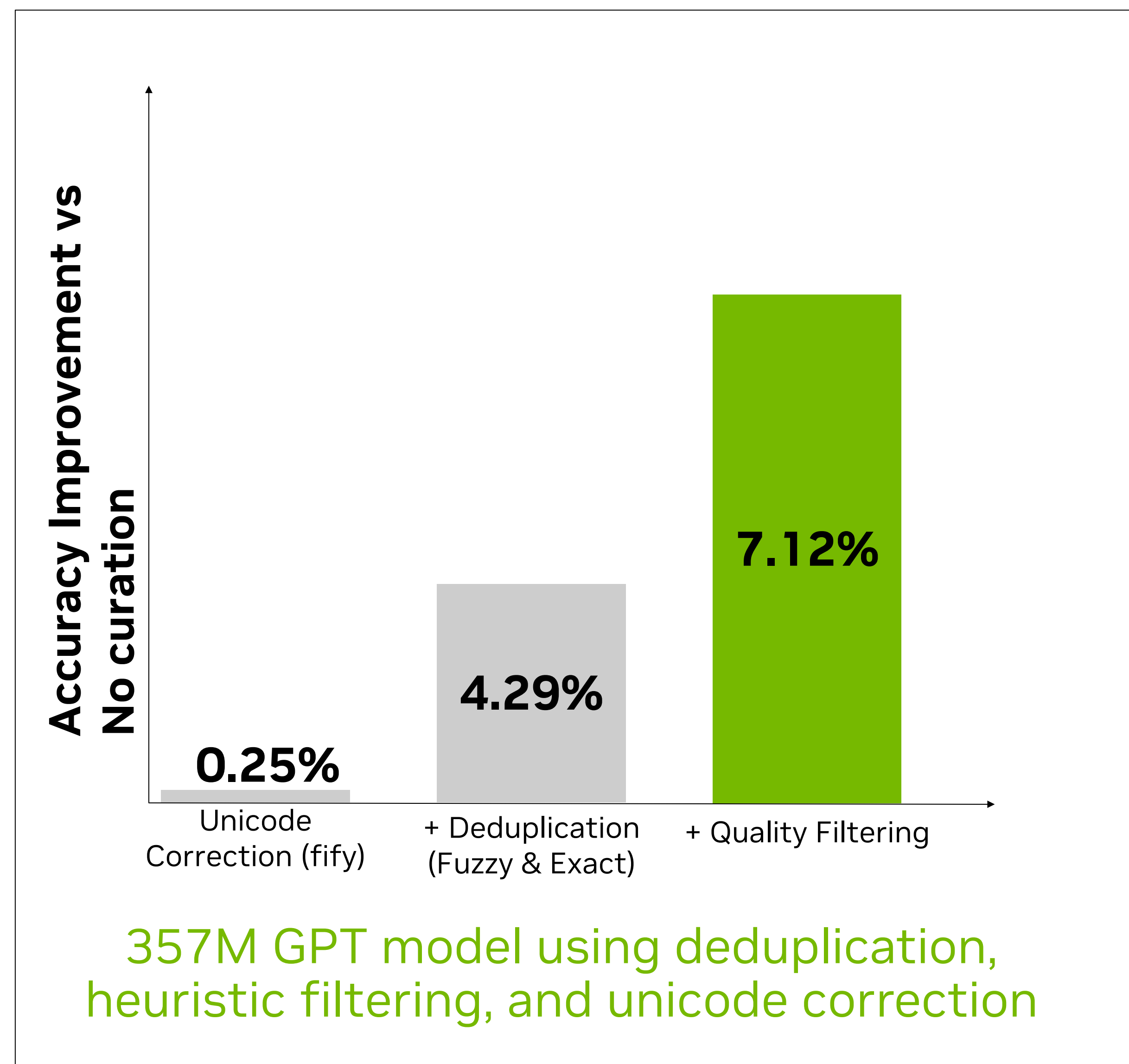| Classification | Mask/remove PII | Acceleration |
|---|---|---|
| **Domain Classifier**<br><br>• Model trained on 1.5 million samples<br><br>• Classifies text in 26 domain classes<br><br><br>**Document Quality Classifier**<br><br>• Classifier quality of document in 'High', 'Medium', 'Low' categories<br><br>• Enables document quality check | • Remove sensitive information from data<br><br><br>• Utilizes State-Of-The-Art spacy model to remove PII information | |

**Acceleration chart:**

Accuracy Improvement vs No curation

- Unicode Correction (fify): 0.25%
- + Deduplication (Fuzzy & Exact): 4.29%
- + Quality Filtering: 7.12%

357M GPT model using deduplication, heuristic filtering, and unicode correction

# Accelerated Data Processing Maximizes LLM Performance & Scale

**Nemo Curator: Up to 7% better Accuracy for LLM Downstream Tasks**



Accuracy Improvement vs No curation

0.25% — Unicode Correction (fify)

4.29% — + Deduplication (Fuzzy & Exact)

7.12% — + Quality Filtering

357M GPT model using deduplication, heuristic filtering, and unicode correction

**Scale to 100+ TB of data**



Processing Time (hrs)

3.4 | 1.7 | 0.9 | 0.5

H100 nodes 80GB: 1 | 2 | 4 | 8

Scaling Fuzzy Deduping 1.2T tokens (RedPajama-v1)

NVIDIA

# NeMo Curator: Example

```python
# Download your dataset
dataset = download_common_crawl("/datasets/common_crawl/", "2021-04", "2021-10", url_limit=10)
# Build your pipeline
curation_pipeline = Sequential([
  # Fix unicode
  Modify(UnicodeReformatter()),
  # Discard short records
  ScoreFilter(WordCountFilter(min_words=80)),
  # Discard low-quality records
  ScoreFilter(FastTextQualityFilter(model_path="model.bin")),
  # Discard records from the evaluation metrics to prevent test set leakage.
  TaskDecontamination([Winogrande(), Squad(), TriviaQA()])
])
# Execute the pipeline on your dataset
curated_dataset = curation_pipeline(dataset)
```

| GitHub | NeMo framework container | PyPI |
| --- | --- | --- |

NVIDIA.

# NeMo Curator - Resources

## Getting Started

- NeMo Framework Container
- GitHub
- PyPI
- Installation Guide
  Developer Page
- User Guide/ Docs
- API Docs
- Classifier Models
- Examples
- Best Practices
- Bugs
- Discussions

## Tutorials & Blogs

### Pre-training / DAPT

- Curating data for LLM training (Blog)
- Curating non-English data (Blog)
- How-to run classifier models
- All blogs

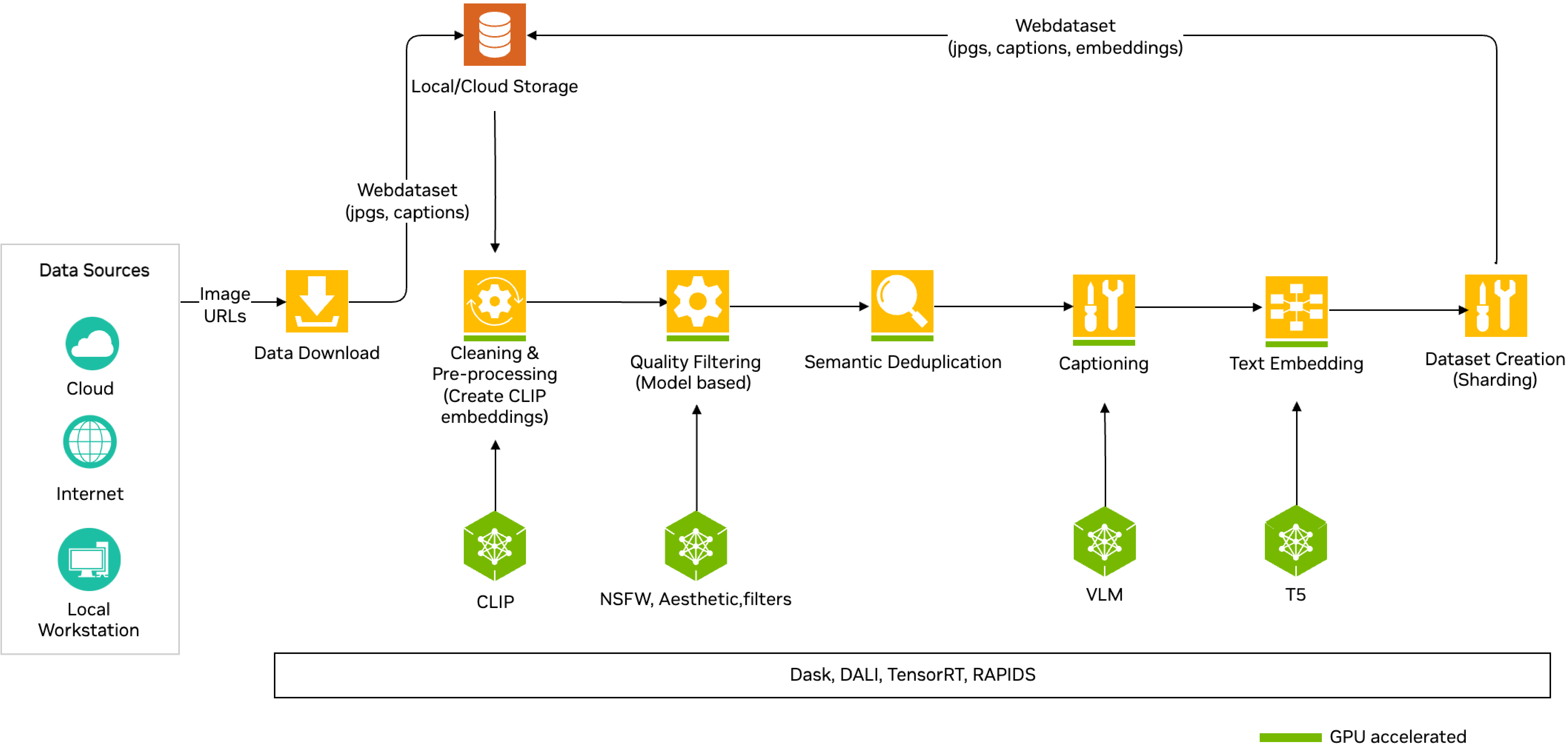### Fine-tuning

- Curating data for PEFT (Blog)
- Curating synthetic data for PEFT
- SDG using Llama 3.1 405B & Nemotron 4-340B
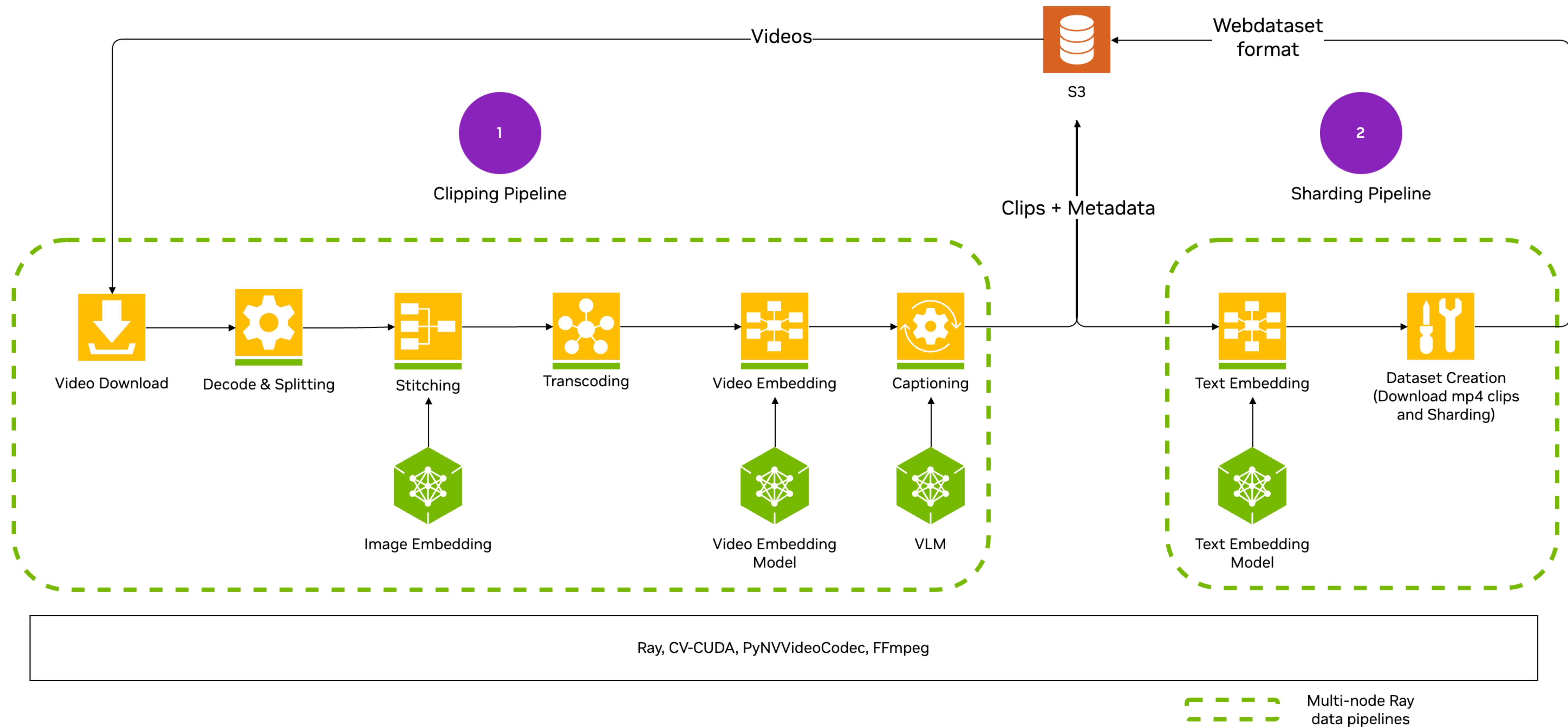- SDG using Nemotron 4-340B

# Image Processing

# NeMo Curator Architecture: Image Processing

# Video Processing

# Video Curation: EA Features



S3

Videos

Webdataset format

Clips + Metadata

**1** Clipping Pipeline

**2** Sharding Pipeline

Video Download → Decode & Splitting → Stitching → Transcoding → Video Embedding → Captioning

Image Embedding

Video Embedding Model

VLM

Text Embedding → Dataset Creation (Download mp4 clips and Sharding)

Text Embedding Model

Ray, CV-CUDA, PyNVVideoCodec, FFmpeg

Multi-node Ray data pipelines

# State of the Art Models

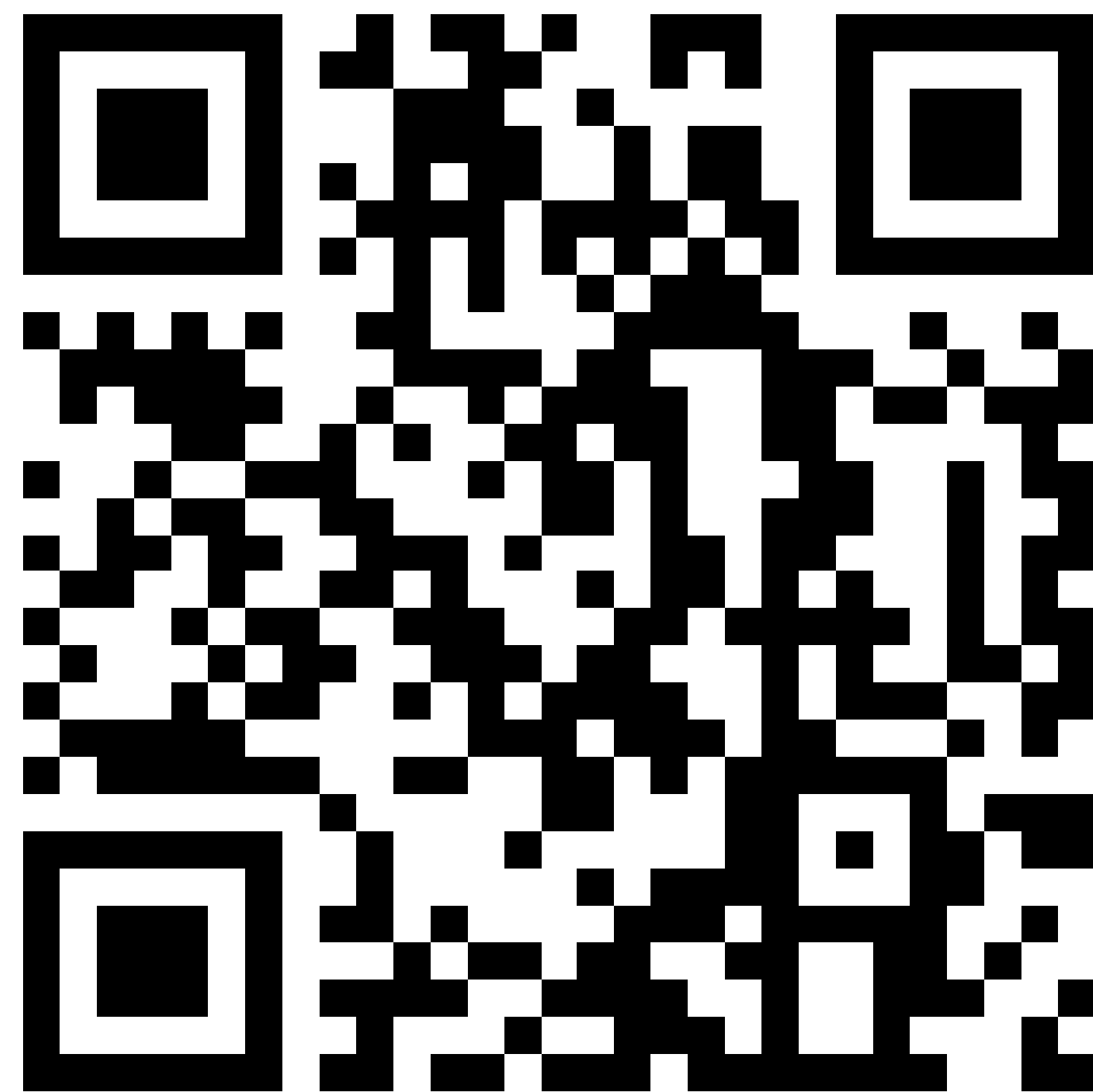Example Use case: Generate new lighting, weather, and geolocations

# Best Practices

- Choosing the Right Quality Model Type

- Handling GPU Out-of-Memory (OOM) Errors

  - Controlling Partition Sizes

- Fuzzy Deduplication Guidelines

  - Reduce bucket counts

  - Reduce buckets per shuffle

  - Adjust files per partition

- GPU Memory and Utilization – Dask Dashboard

# Developer Tools and Resources

## Accelerate innovation and growth

Learn more: developer.nvidia.com

## Individuals

### Software
100s of APIs, models, SDKs, microservices, and early access to NVIDIA tech

### Training
Hands-on self-paced courses, instructor-led workshops, and certifications

### GPU Sandbox
Approval basis, multi-GPU and multi-node

### Learning
Tutorials, self-paced courses, blogs, documentation, code samples

### Community
Dedicated developer forums, meetups, hackathons

### Ecosystem
GTC, NVIDIA Partner Network

## Organizations

### Startups
Cloud credits, engineering resources, technology discounts, exposure to VCs

### Venture Capital
Deal flow and portfolio support for Venture Capital firms

### Higher Education
Teaching kits, training, curriculum co-development, grants

### ISVs and SIs
Engineering guidance, discounts, marketing opportunities

### Research
Grant programs, collaboration opportunities

### Enterprises
Tailored developer training, skills certification, technical support

NVIDIA.

# Thank You