

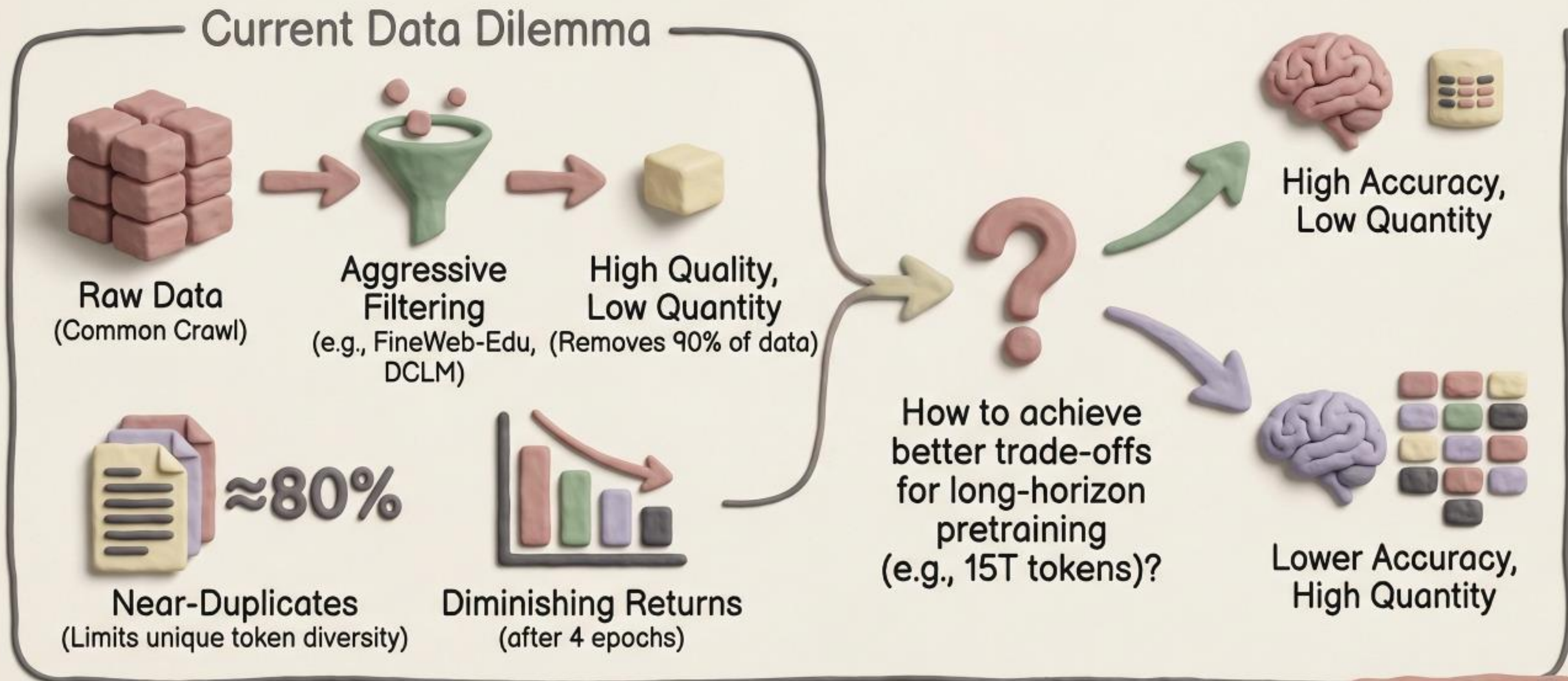
# Nemotron-CC:

Transforming Common Crawl  
into a Refined Long-Horizon  
Pretraining Dataset





# The Challenge: Long-Horizon LLM Training





# Research Contributions & Key Results

## Main Contributions

- 1. 6.3T token dataset transformation method
- 2. Proven effectiveness through comprehensive comparisons
- 3. Detailed ablation studies revealing best practices

## Key Results



15T  
Tokens



1.1T high-quality subset: +5.6 MMLU over DCLM



Full 6.3T dataset matches DCLM, 4x more unique tokens

Nemotron-CC 8B (15T)

+3.1 ARC-Challenge

+0.5 Average (10 tasks)

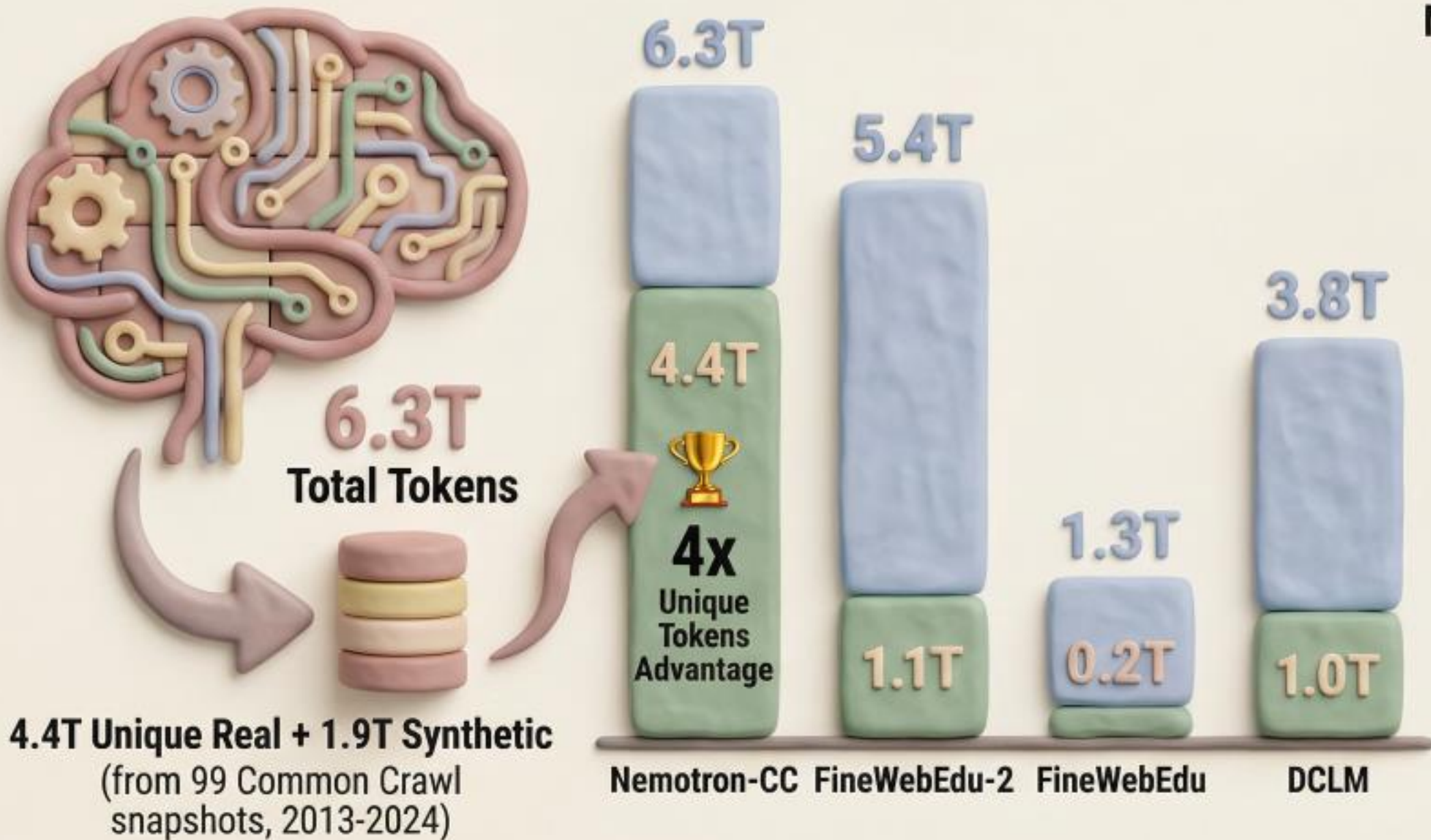
Outperforms Llama 3.1 8B



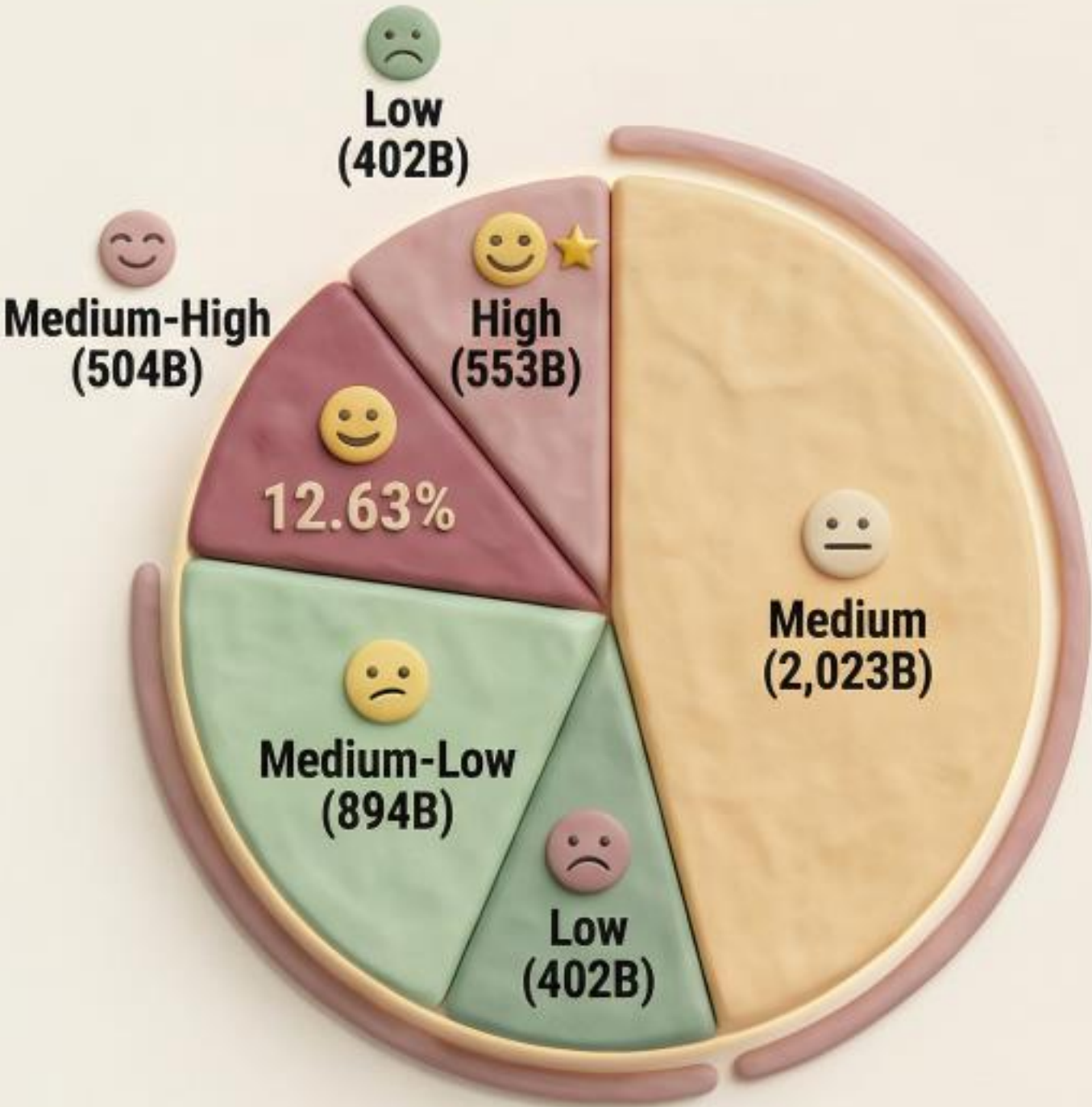
# Dataset Overview: Scale & Composition

Unprecedented Scale with a Focus on Unique, High-Quality Data

## Total Size & Unique Tokens



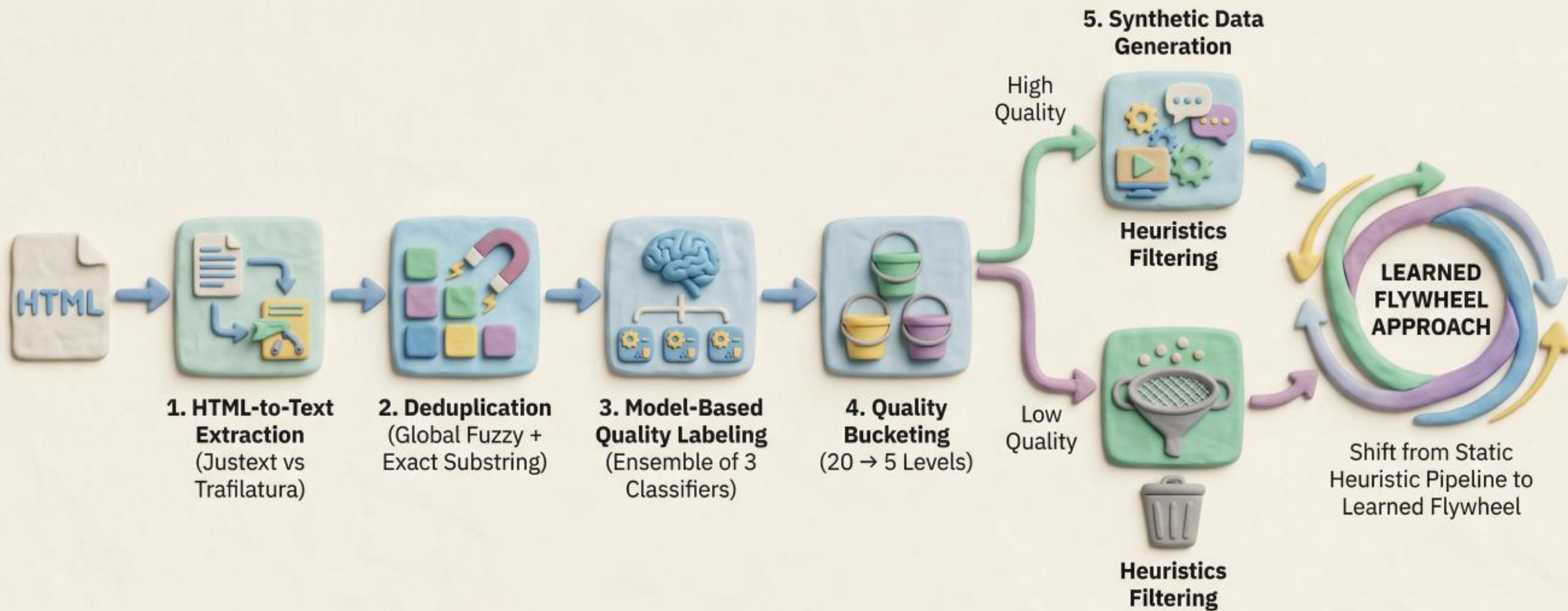
## Quality Distribution



Over 70% of the dataset is Medium quality or higher, ensuring robust training.



# The Nemotron-CC Pipeline: An Overview





# HTML-to-Text Extraction: Maximizing Token Yield

Comparing Trafilatura vs Justext for Pretraining Data



## Trafilatura Extraction



## Justext Extraction (Recommended)



## Key Insight & Explanation

-  **Prioritize Absolute HQ Tokens over Percentage:** Quality bucketing enables exact control during training.
-  **Justext extracts more tokens with similar perceived quality.** Ablation results show no negative impact on downstream accuracy.



NVIDIA Research

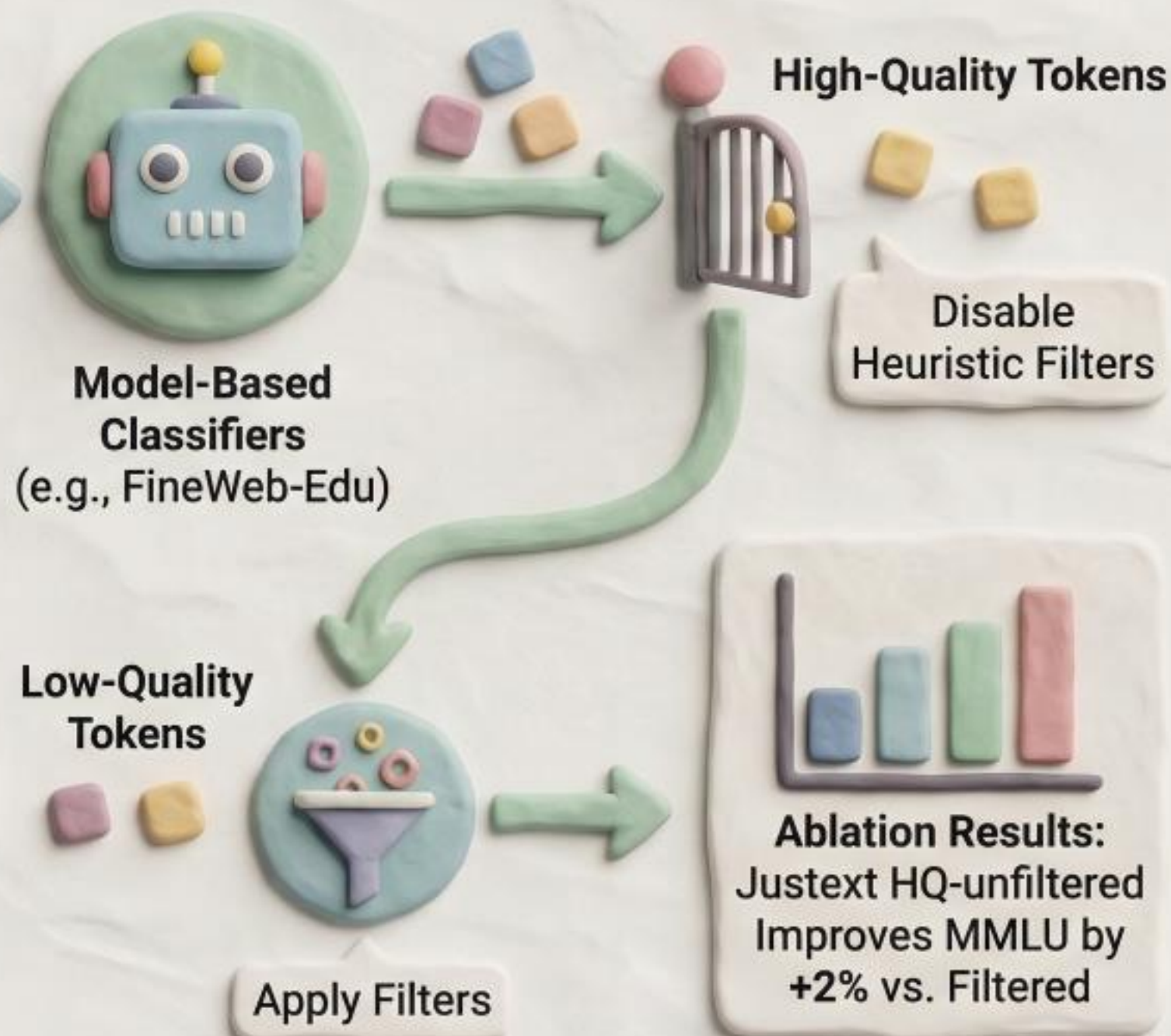


# Heuristic Filtering: A Strategic Reevaluation

## Conventional Approach



## Novel Strategy





# Model-Based Quality Labeling: Classifier Ensemble

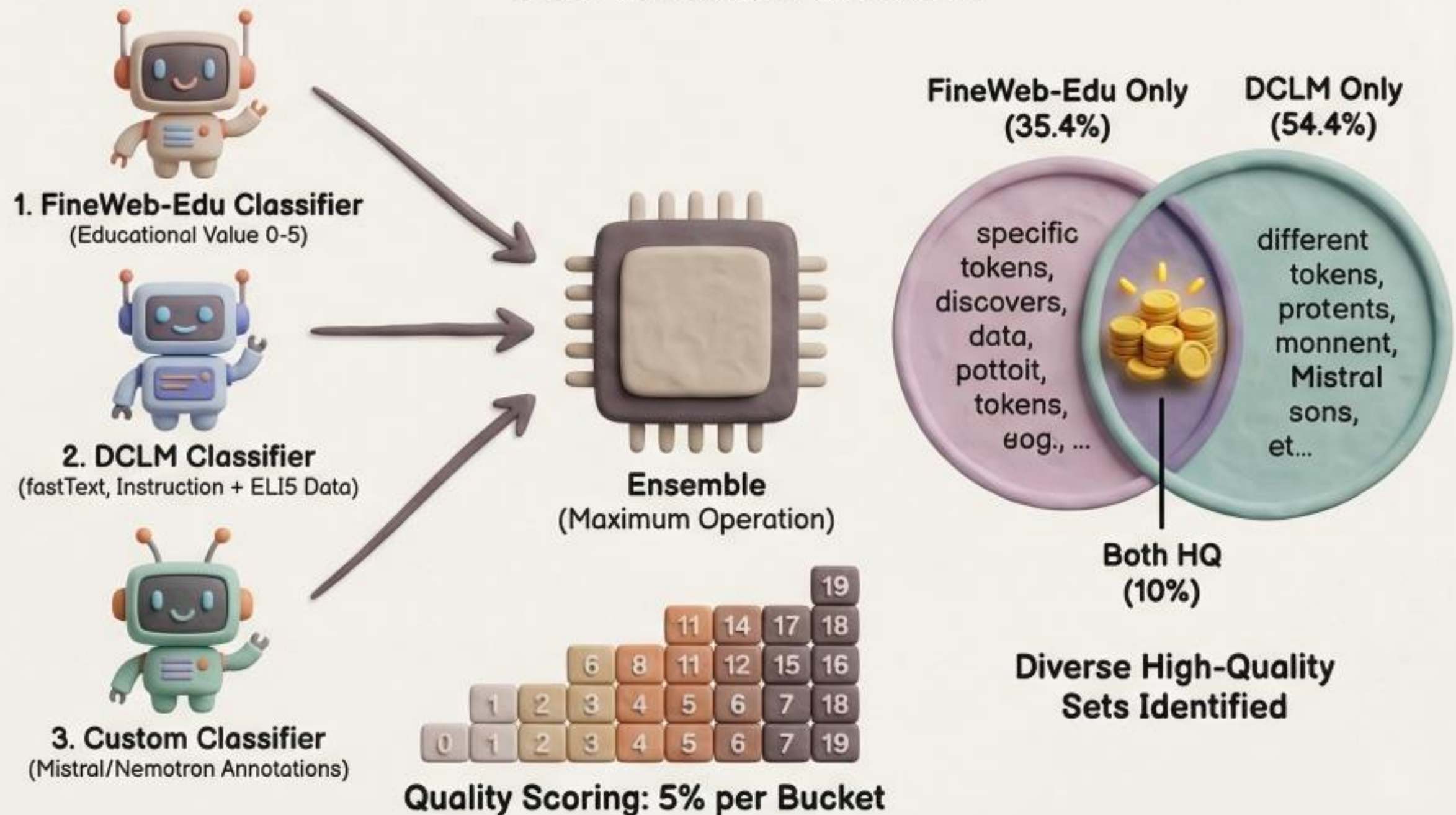
Overcoming Single Classifier Bottlenecks with a Three-Pronged Approach

## The Bottleneck Problem



The Bottleneck Problem

## The Ensemble Solution





# Quality Bucketing: From Scores to Training Labels

20 Fine-Grained Buckets (0-19)



Bucketing Methodology



66% Default



Annealing Evaluation (Regrouping Process)

34% Evaluated Bucket



Final 5-Category Grouping

**High**  
(Bucket 19, 553B, 12.63%)

**Medium-High**  
(Bucket 18, 504B, 11.52%)

**Medium**  
(Buckets 12-17, 2,023B, 46.24%)

**Medium-Low**  
(Buckets 7-11, 894B, 20.43%)

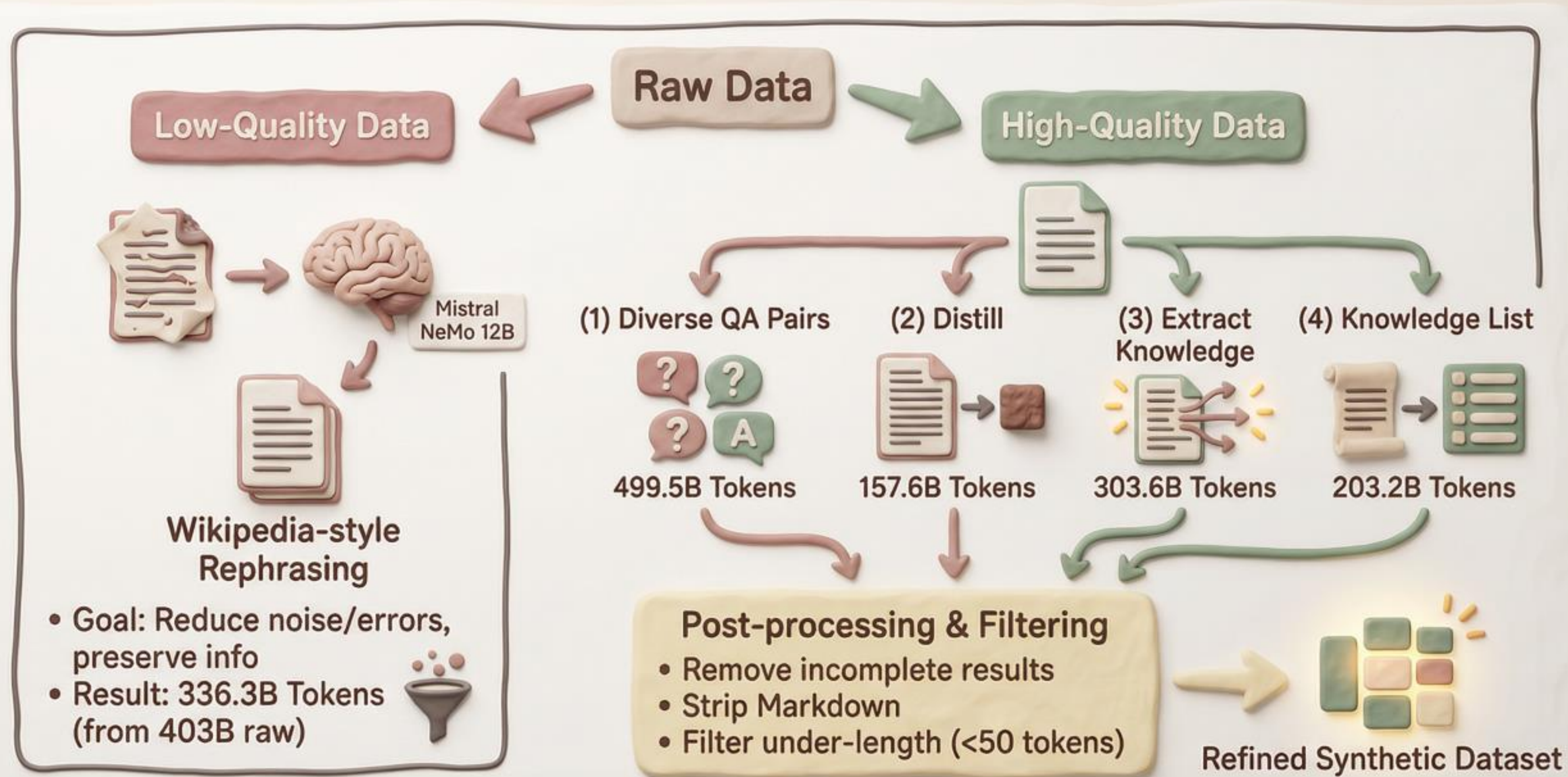
**Low**  
(Buckets 0-6, 402B, 9.18%)



Aligned with Downstream Performance



# Synthetic Data Generation: Dual Strategies





# Synthetic Data Examples & Quality Control

## Generated QA Examples & Types

Question: "Which year did the UN implement the 2030 agenda for SDGs?"

Answer: "January 1, 2016"



Factual Recall



Conceptual Understanding



Multiple-Choice



Yes/No

## Post-Processing Pipeline



Remove Incomplete



Eliminate Markdown (\*\*)



Strip Prefixes



Remove Quotes



Filter <50 Tokens

Concatenate Passages



Wikipedia/QA Handling



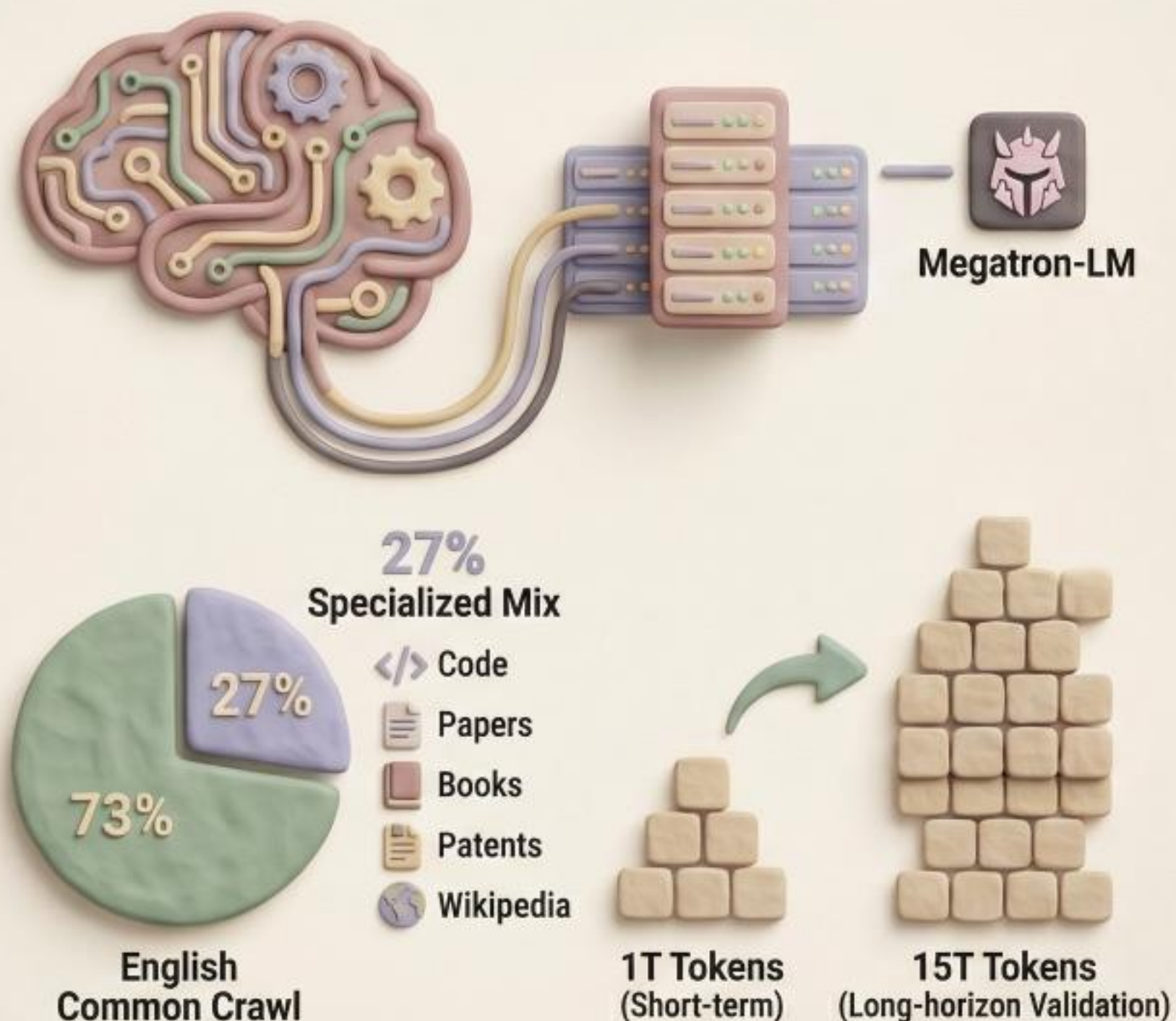
Shuffle & Retain

Append to End

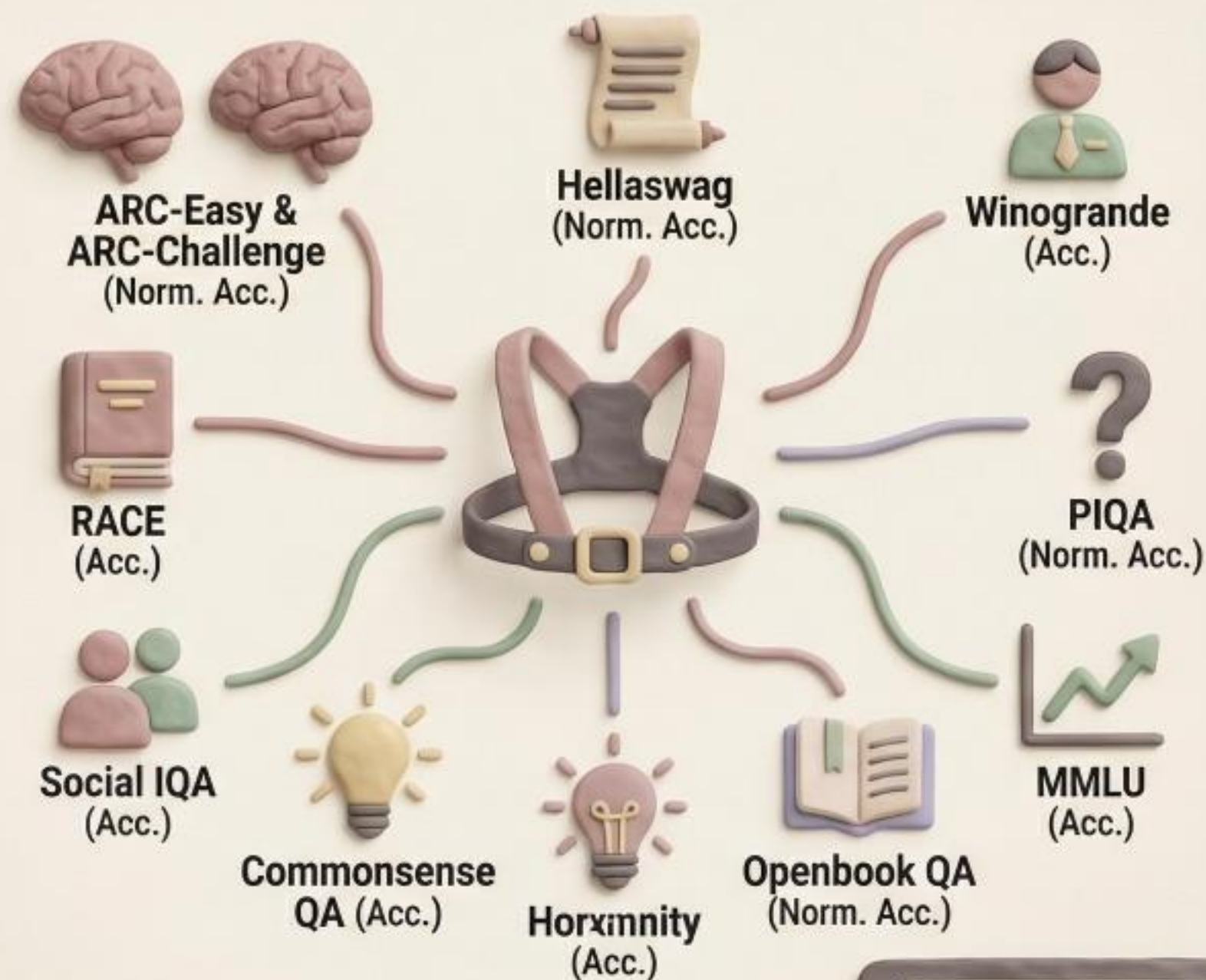


# Experimental Setup: Training & Evaluation

## Training Configuration & Data



## Evaluation Methodology



Hyperparameters detailed in Appendix D.



# Short Token Horizon Results (1T Tokens)



8B Model  
on 1T Tokens

Nemotron-CC-HQ



DCLM



Classifier  
Ensembling &  
Synthetic Data



Nemotron-CC  
(full dataset)

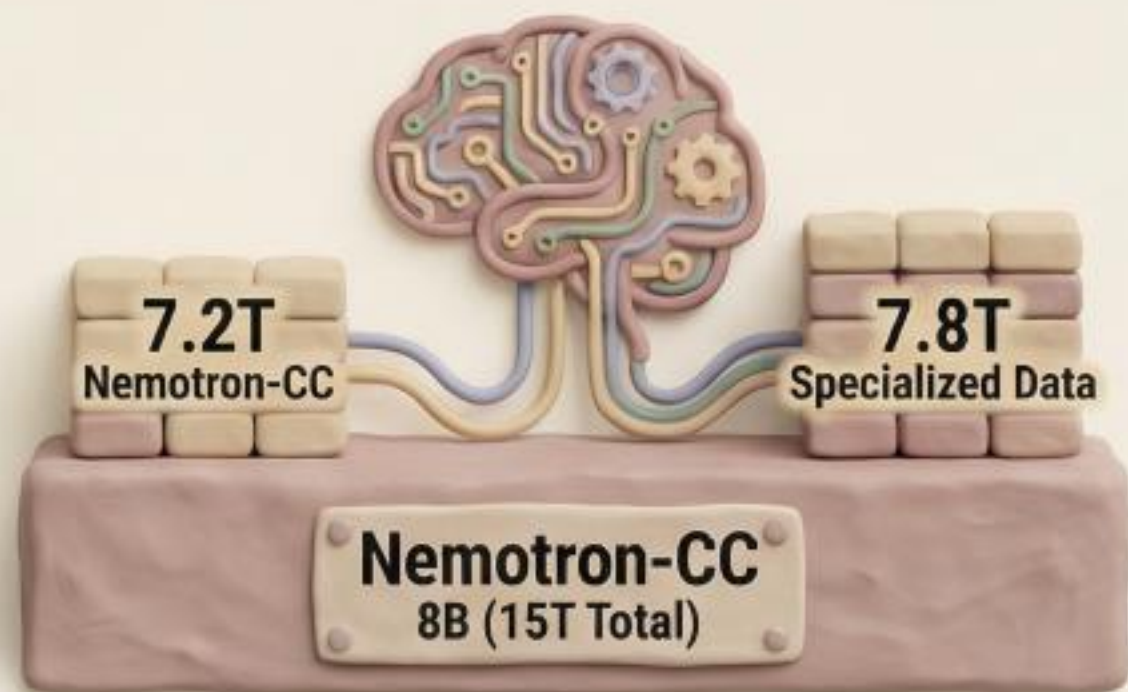


DCLM  
(baseline)

Key Insight: Classifier ensembling and synthetic data effective even in non-data-constrained settings. Superiority demonstrated across most tasks.



# Long Token Horizon Results (15T Tokens)

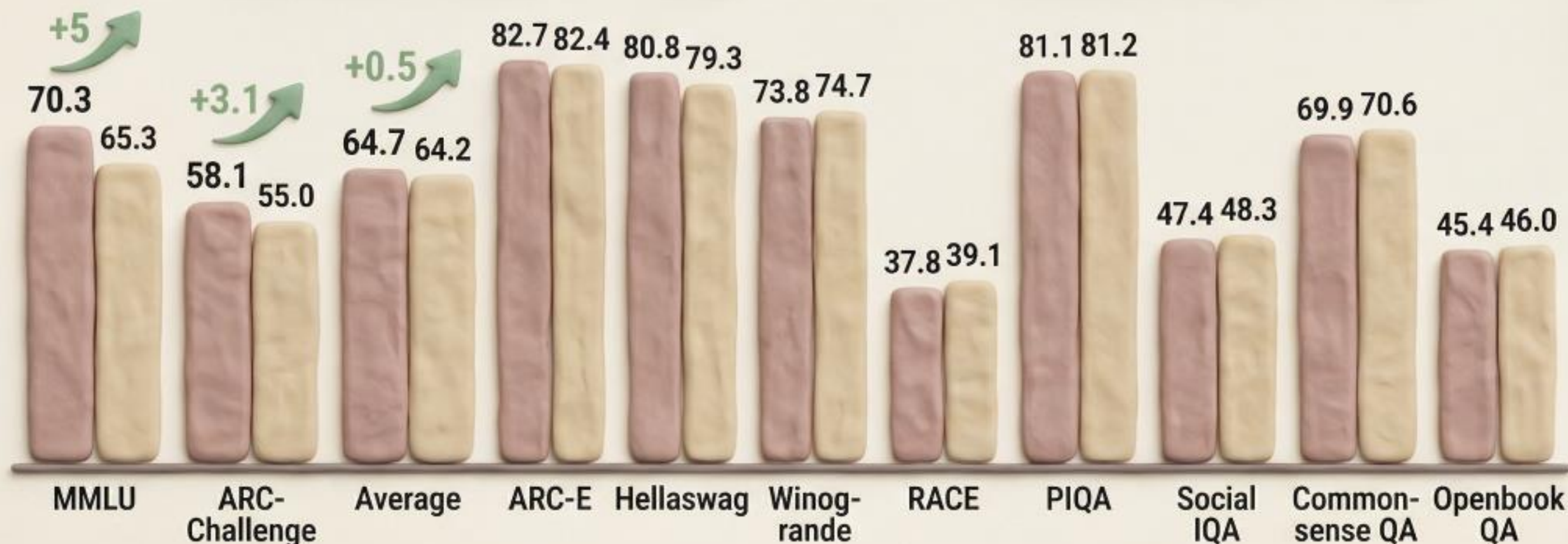


## Key Achievements

**+5** MMLU

**+3.1** ARC-Challenge

**+0.5** Average



**Validates dataset suitability for state-of-the-art long-horizon training.**

Note: Evaluation uses lm-evaluation-harness; Meta's numbers may differ due to customizations.



# Ablation Study: Extractor & Filter Impact

## Configurations & Performance

Trafilatura filtered	MMLU: 55.4	Average: 60.6	✓
Justext filtered	MMLU: 54.1	Average: 60.9	🔍
Justext unfiltered	MMLU: 55.5	Average: 60.3	📦
Justext HQ-unfiltered	MMLU: 57.5	Average: 60.6	★

Diagram illustrating the configurations and performance metrics (MMLU and Average) for different extractor and filter combinations. The configurations are ranked from top to bottom: Trafilatura filtered, Justext filtered, Justext unfiltered, and Justext HQ-unfiltered. The Justext HQ-unfiltered configuration is marked as the 'Best' configuration with a star and an upward arrow.

Annotations: Finn (between Trafilatura filtered and Justext filtered), Com (between Justext filtered and Justext unfiltered), Best (between Justext unfiltered and Justext HQ-unfiltered).

## Key Findings & Strategic Insight

- Justext yields 57.4% more HQ tokens than Trafilatura without accuracy impact



- Removing filters from HQ tokens improves MMLU by +2%



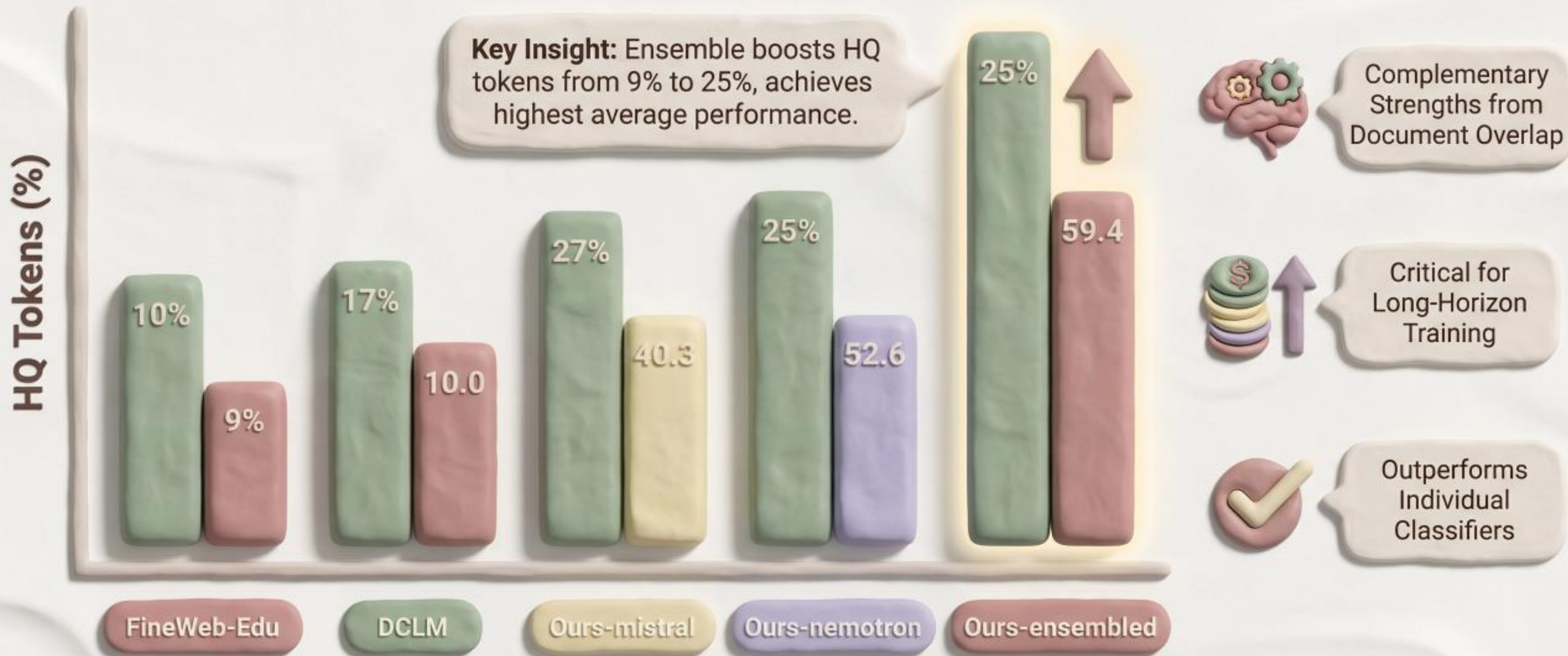
- Combining Justext + no HQ filtering maximizes both token yield and accuracy



**Strategic Insight:** Filters remove 18.1% of HQ tokens unnecessarily.  
**Conclusion:** Apply heuristics only to low-quality tokens.



# Ablation Study: Classifier Comparison & Ensemble





# Ablation Study: Synthetic Data Evaluation

## Low Quality (LQ) Data Rephrasing



**LQ-Base**  
(Original CC)

MMLU: 48.2 | Avg: 52.5

**LQ-Synthetic**  
(Rephrased)

MMLU: 47.1 | Avg: 54.0

**+1.50**  
Average  
Score Boost

↑ ARC-E: +3.6  
↑ OBQA: +3.6  
↑ CSQA: +4.7



Potential  
Misinformation  
Risk

## High Quality (HQ) Data Augmentation



**HQ-Base**  
(8x HQ)



MMLU: 53.4 | Avg: 55.8

**HQ-Synthetic**  
(4x HQ + Synthetic)

MMLU: 53.6 | Avg: 56.7

**+0.9**  
Average  
Score Boost

📄 Fresh Unique Tokens  
📈 Diverse Styles for QA  
📊 Outperforms 8 HQ Epochs

 **Key Findings:** Synthetic data improves average scores by providing fresh tokens and diverse styles, especially for specific abilities like QA. However, data curation is crucial to mitigate noisy examples and potential accuracy drops. 

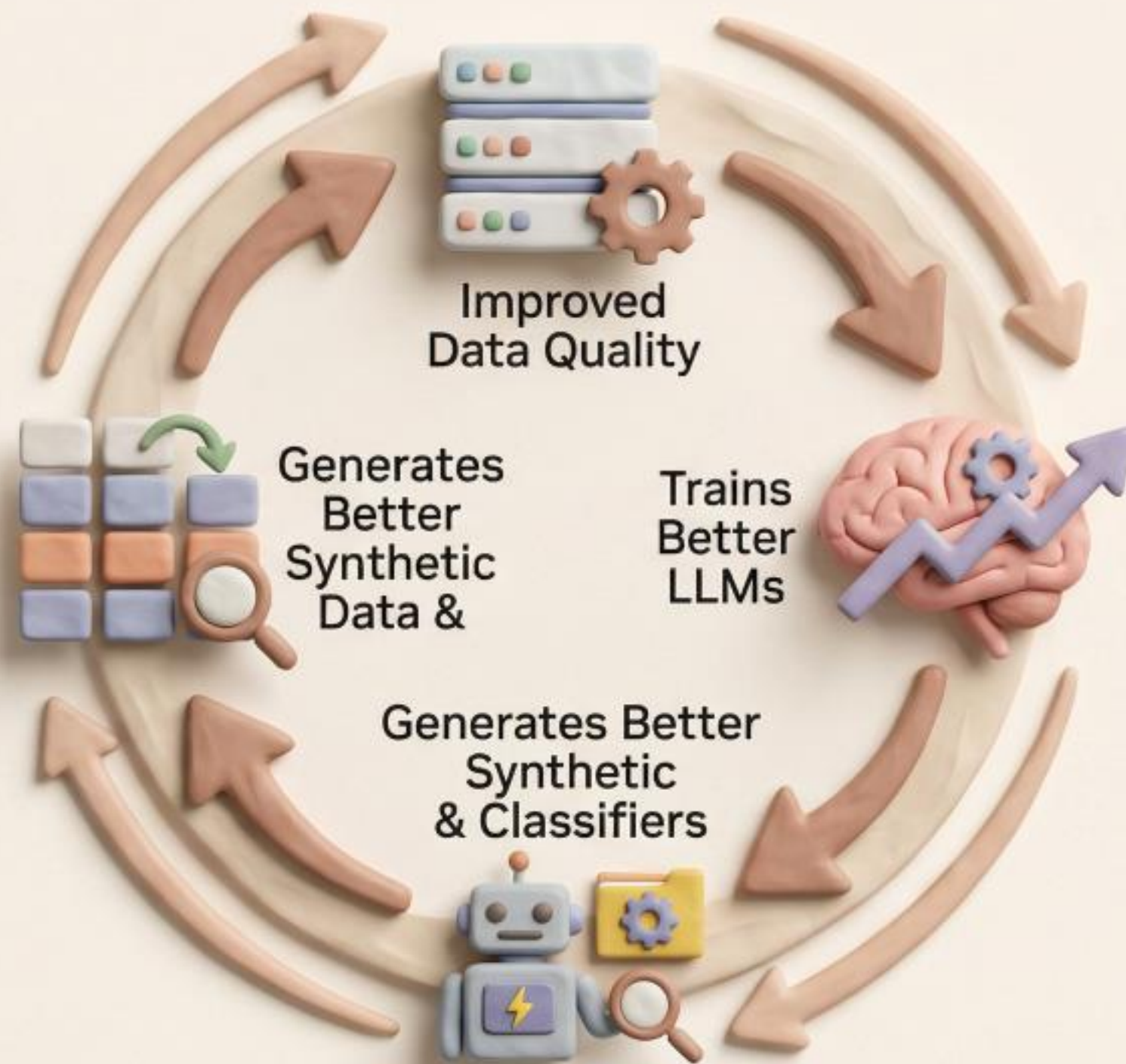


# Guiding Principle: The Learned Flywheel

Traditional Approaches  
(Static Heuristics)



Traditional Approaches  
(Static Heuristics)



Dynamic &  
Self-Improving



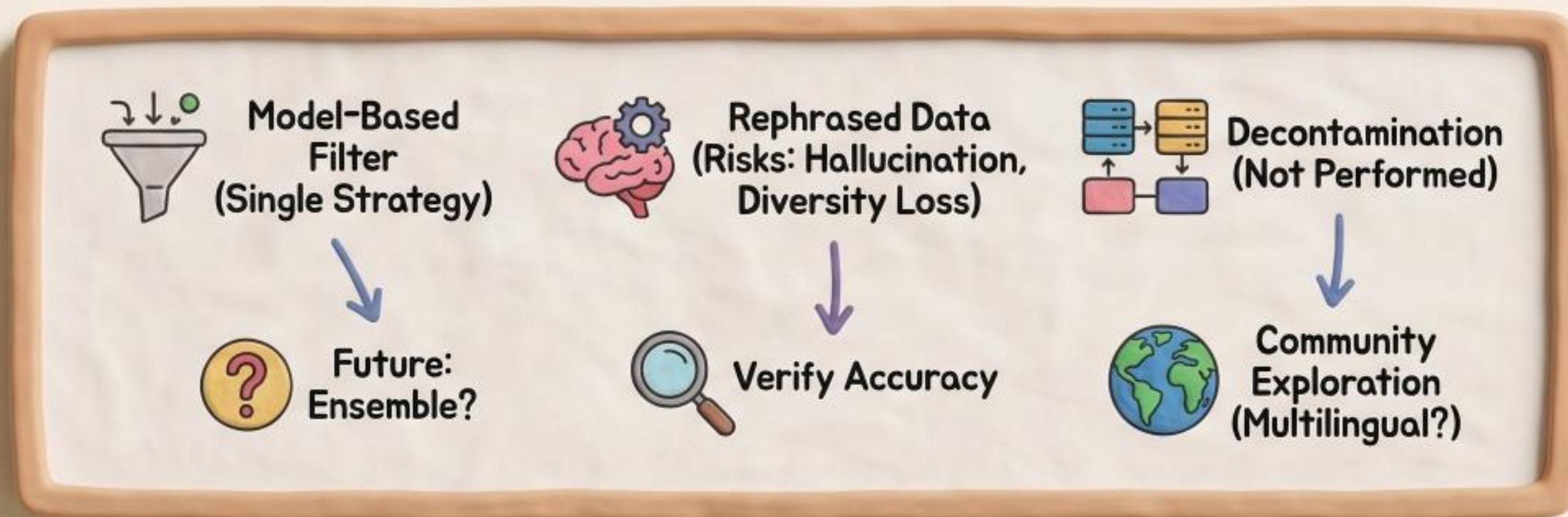
Future-Proof:  
Accuracy Naturally  
Increases



Reduces Manual  
Intervention



# Limitations & Future Directions



## Filtering & Pipeline

- Single filter strategy tested; needs ensemble for higher quality.
- Pipeline components (e.g., lang ID) not fully ablated.

## Rephrasing & Scope

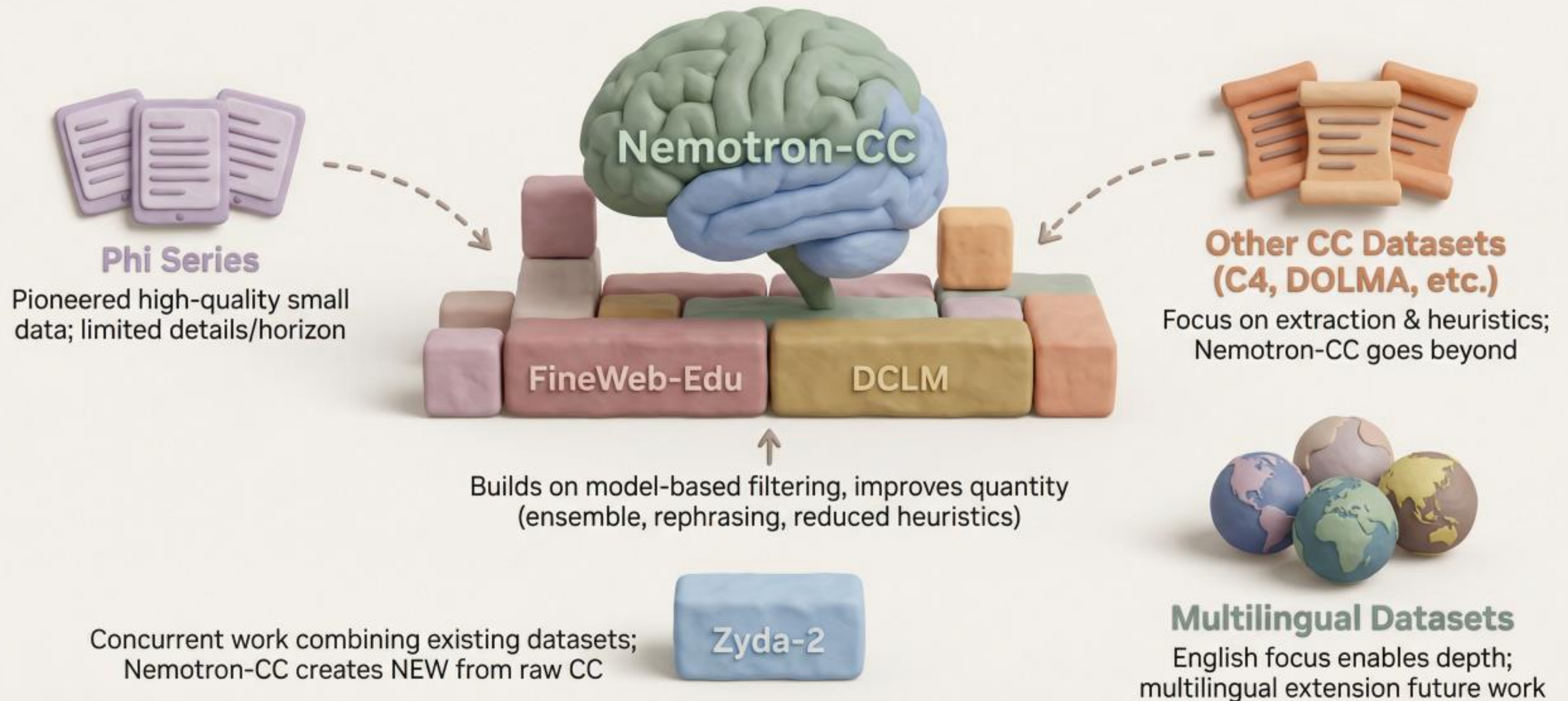
- Factual accuracy & fidelity of rephrased data unverified.
- Tested on English only; multilingual extension needed.

## Decontamination

- Dataset not decontaminated due to lack of consensus.
- Requires community exploration on best practices.



# Related Work: Building on Previous Research





# Synthetic Data Literature & Inspirations



## Instruction Pre-training (Cheng et al. 2024)

Synthesized instruction-response pairs for pre-training.



## TinyStories (Eldan and Li 2023)

Small models trained on synthetic short stories generate fluent narratives.



## Textbook Models (Gunasekar et al. 2023)

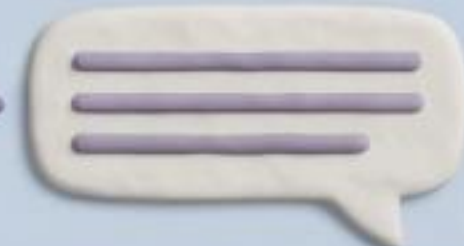
Synthetic textbooks and exercises achieve impressive coding benchmarks.

## Rephrasing Approach (Maini et al. 2024)



Qwen

jumbled to : /a; jum-  
jum, text, meth, but  
task is to ewen in ...



Mistral

Smaller models (Qwen-1.8B, Mistral-7B)  
adequate for web data rephrasing.

## Nemotron-CC Extension



Adopts and extends: more prompts, specialized for quality tiers, large-scale demonstration



# Dataset Availability & Open Source Release

## Public Release & Access



[data.commoncrawl.org/contrib/Nemotron...](https://data.commoncrawl.org/contrib/Nemotron-4-edu-classifier)

Released under Common Crawl Terms of Use.

## Open Source Tools



[github.com/NVIDIA/NeMo-Curator](https://github.com/NVIDIA/NeMo-Curator)

Reference implementation:  
Apache 2.0 NeMo Curator library.

## Quality Classifier Models

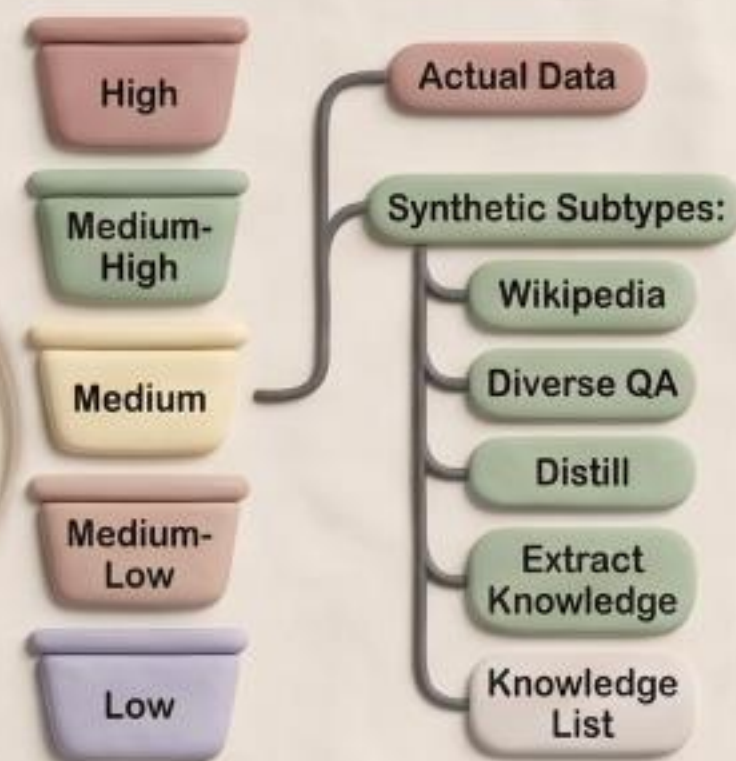


[nemocurator-fineweb-nemotron-4-edu-classifier](#)

[nemocurator-fineweb-mixtral-edu-classifier](#)

Released on HuggingFace.

## Dataset Organization



Enables community experiments:  
Quality vs Diversity & Curriculum Design.



# Training Curriculum: Two-Phase Approach

## Phase 1: 9T Tokens

Medium  
Medium-High  
High Quality  
(real + synthetic)  
English Common Crawl  
(5.31T - 59%)



**59% English Common Crawl (5.31T).**  
Uses medium, medium-high, and high quality data  
(real + synthetic).

## Phase 2: 6T Tokens

High Quality  
(real + synthetic)  
English Common Crawl  
(1.86T - 31%)



**6T tokens, 31% English Common Crawl (1.86T)**  
Uses only high quality data (real + synthetic).

**Key Insight:**  
4-8 epochs of high-quality  
data optimal before medium-  
quality benefits outweigh.

**Combined Total:**  
**47.8% English CC (7.17T)**

### Non-CC Portion (27%)

Books & Patents (9%)	Papers (9%)	Code (5%)	Conversational (3%)	Wikipedia (1%)



# Technical Implementation: Tools & Infrastructure

## HTML & Language ID



**Justext**  
(Pomikálek 2011)



**Trafilatura**  
(Barbaresi 2021)



**Language ID:**  
pyclid2  
& FastText  
(lid176, thresh 0.3)



## Deduplication



**NeMo Curator**  
(Global Fuzzy)

**deduplicate-text-datasets**  
(Exact Substring)

## Quality Classifiers

**FineWeb-Edu**



**DCLM**

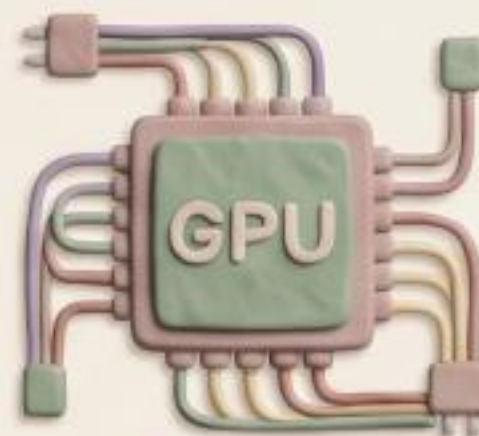


**Snowflake-arctic-embed-m**  
(Custom)



**'quality check'**

## Training & Generation



**Megatron-LM**  
(Training)

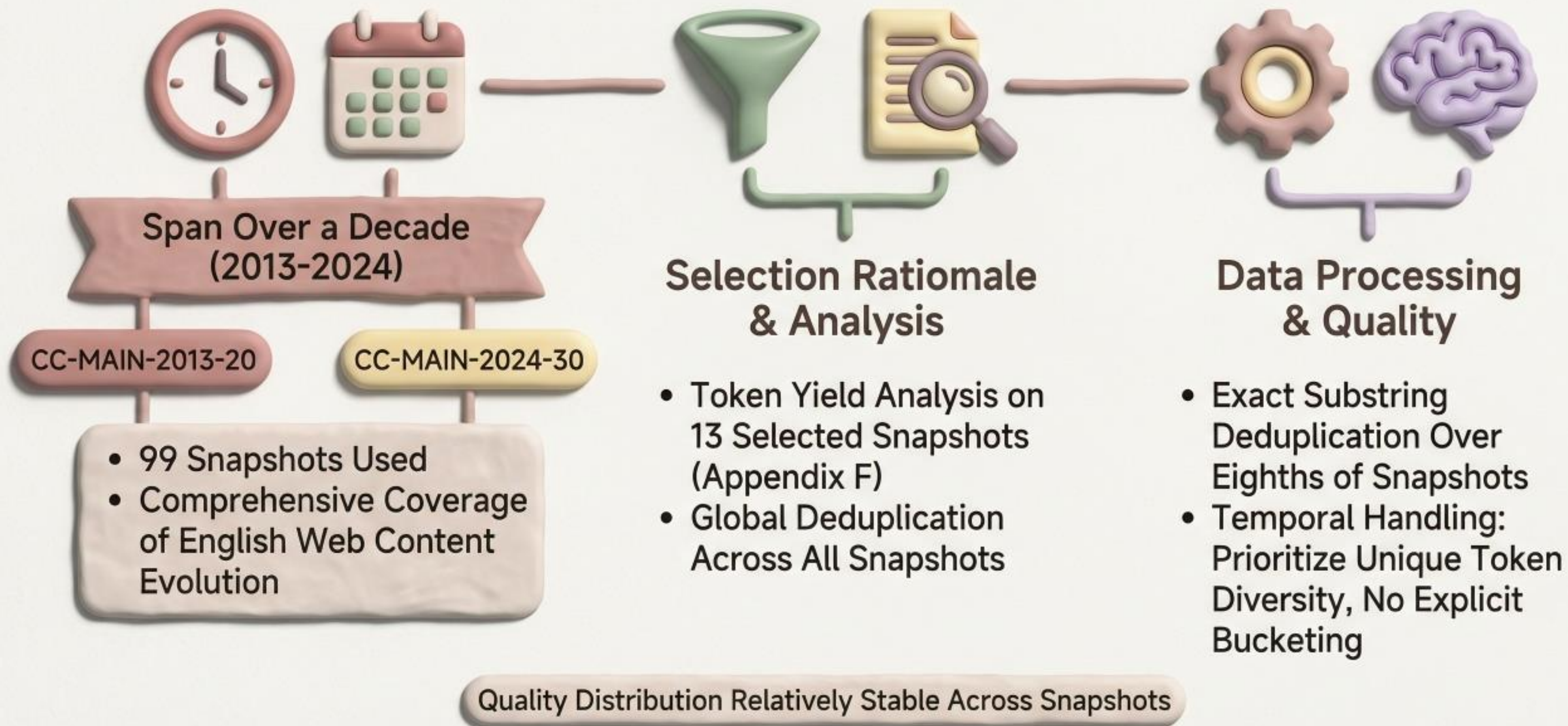


**Mistral NeMo 12B**  
(Synthesis, TensorRT-LLM, NeMo-Skills)

**LM**  
Evaluation  
Harness



# Common Crawl Snapshot Selection





# Quality Control: Deduplication Strategy



## Comprehensive Approach

### 1 Global Exact



- Remove identical docs across 99 snapshots

### 2 Global Fuzzy



- Near-duplicate detection (similarity hashing).  
Tool: NeMo Curator

### 3 Exact Substring



- Split chunks into 8 parts, deduplicate within.  
Tool: deduplicate-text-datasets



## Rationale & Impact

- Diminishing returns after 4 epochs (Muennighoff et al., 2024)



Enables longer effective training without repetition penalty



Trade-off: Some near-duplicates may preserve style



# Evaluation Metrics: Comprehensive Benchmark Suite

## Reasoning & Commonsense



- ARC-Easy & ARC-Challenge
- Hellaswag
- Winogrande
- PIQA
- Social IQA
- Commonsense QA
- Openbook QA

## Reading Comprehension



- RACE (middle/high school exams)

## Knowledge & Reasoning



- MMLU (57 subjects across STEM, humanities, social sciences)

## Metrics & Coverage



Normalized &  
Raw Accuracy

Science, Physical,  
Social, Reading,  
Broad Knowledge

## Standardized Evaluation



Im-evaluation-harness for fair comparison  
Reported numbers may differ due to implementation.



# Key Takeaways: What Makes Nemotron-CC Work

## Classifier Ensembling



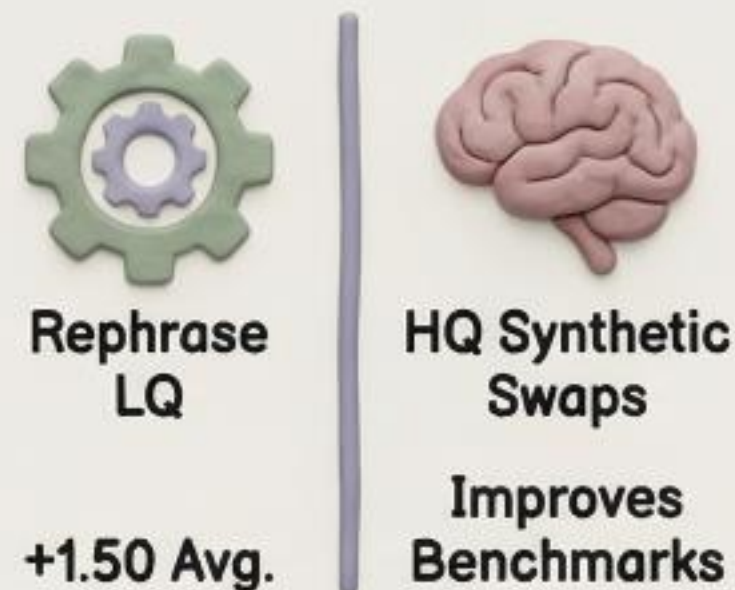
## Strategic Filter Removal



## Justext Extraction



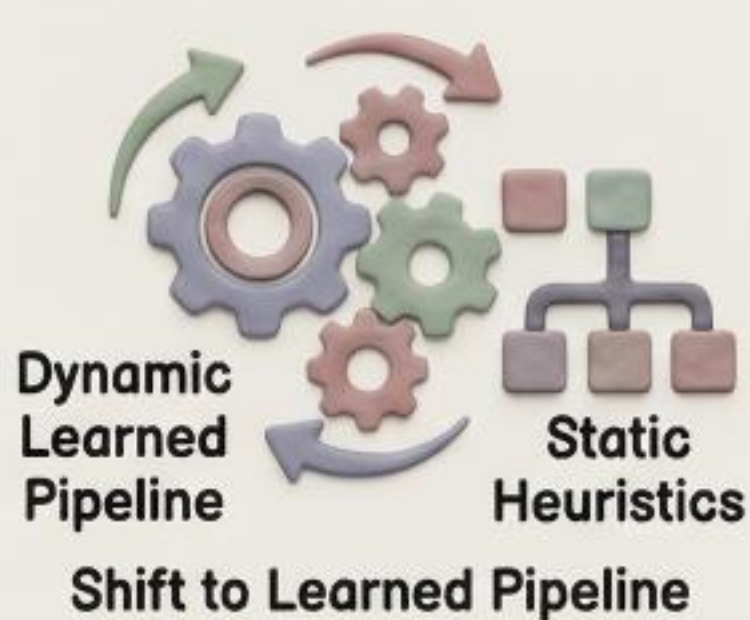
## Synthetic Data Strategies



## Long-Horizon Focus



## Learned Flywheel



## Curriculum Flexibility



## Overall Impact





# Impact on the LLM Community



## RESEARCH & INDUSTRY IMPLICATIONS



## DEMOCRATIZING PRETRAINING

- Public 6.3T Token Dataset (Proven Quality)
- Enabling Smaller Labs: Competitive Models without Massive Proprietary Data
- 15T+ Token Training Feasible for 8B+ Models



## NEW QUALITY-QUANTITY BALANCE

- Challenging Aggressive Filtering Assumption
- Quality-Aware Synthesis as Reusable Technique



## METHODOLOGICAL CONTRIBUTIONS & OPEN SCIENCE

- Classifier Ensembling, Strategic Filtering
- Full Pipeline, Classifiers, Ablations Released for Reproducibility
- Foundation for Future Research: Quality-Bucketed Structure



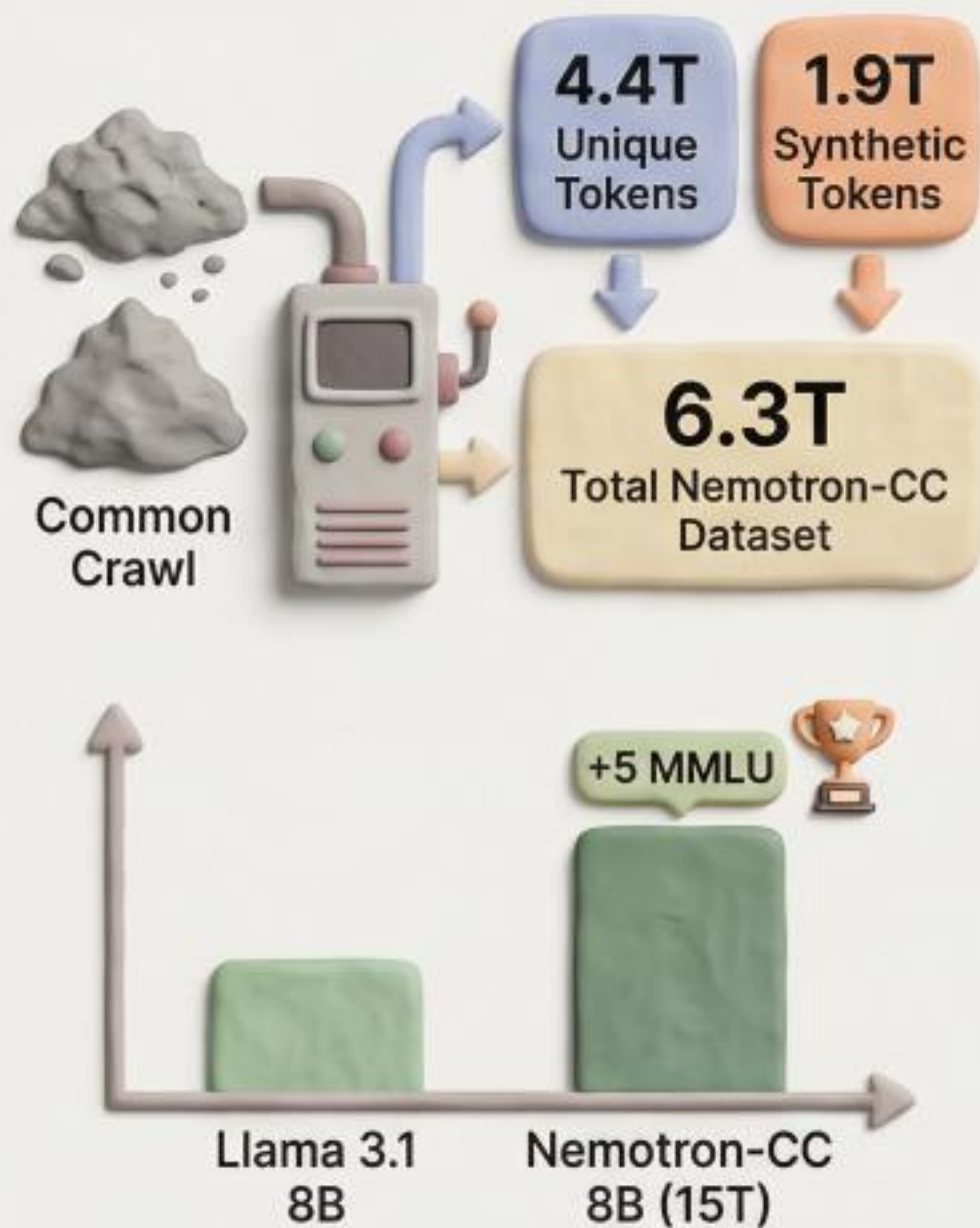
## INDUSTRY IMPACT & INNOVATION

- Better Open Models Drive Competition
- Long-Horizon Training Unlocked

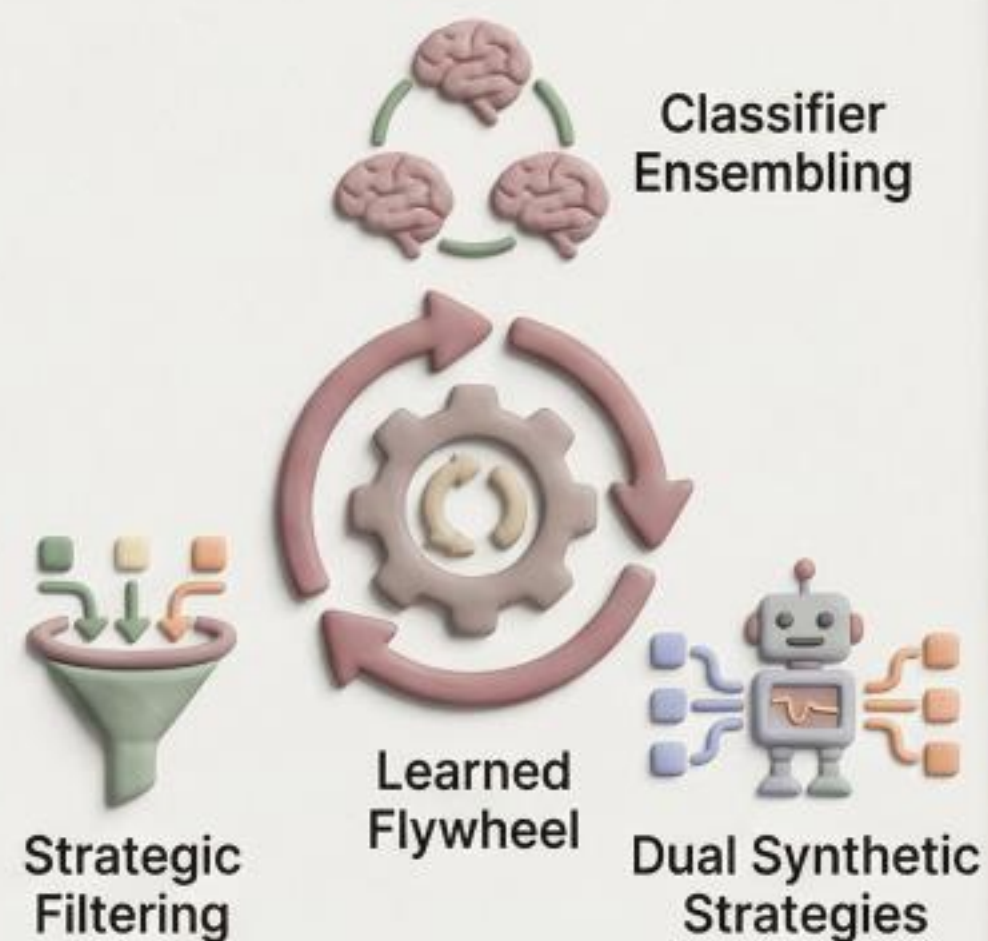


# Conclusion: Advancing Long-Horizon Pretraining

## Key Achievement & Result



## Core Innovations & Principle



Shift from static pipelines to continuous-improvement learned systems.

## Impact & Future Direction



- ➡ Open Access Granular Dataset
- ➡ Enables State-of-the-Art Long-Horizon Training
- ➡ **Future:** Multilingual Extension
- ➡ **Future:** Improved Ensembling
- ➡ **Future:** Factual Verification



**Future Work**