# DEMYSTIFYING SYNTHETIC DATA IN LLM PRE-TRAINING

## A Systematic Study of Scaling Laws, Benefits, and Pitfalls

# Synthetic Data Generation Paradigms

## 1. Web Rephrasing (HQ/QA)

**High-Quality (HQ)**
Wikipedia-style clarity

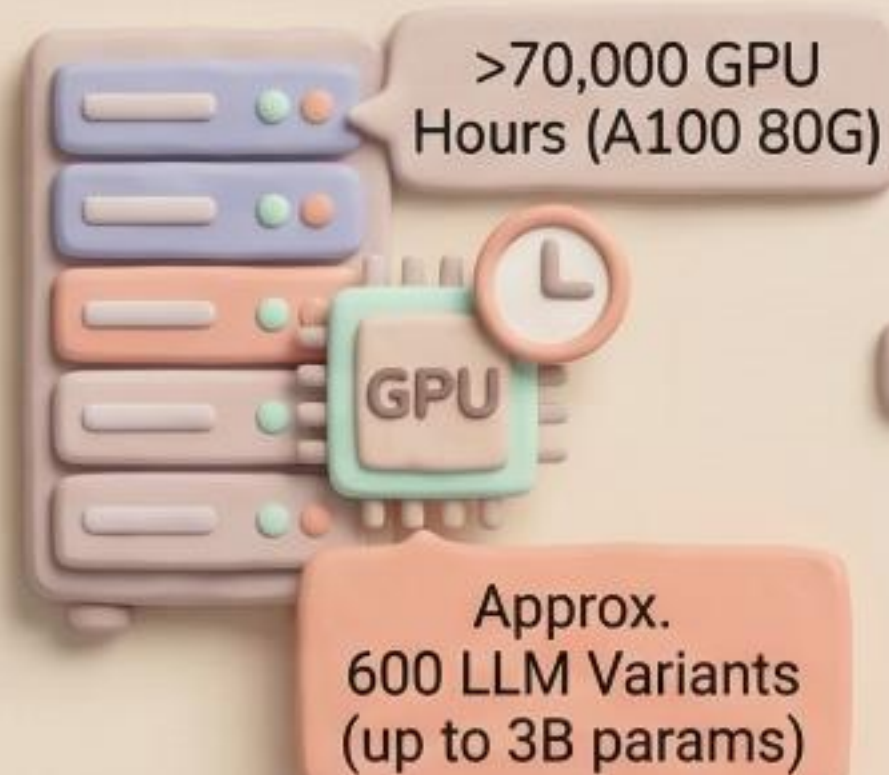**Question-Answering (QA)**
Conversational formats

## 2. Synthetic Textbooks (TXBK)

- Novel educational content
- Dense information
- Examples & exercises

# Experimental Setup & Scale

## Massive Infrastructure

>70,000 GPU Hours (A100 80G)

GPU

Approx. 600 LLM Variants (up to 3B params)

- >70,000 GPU Hours (A100 80G)
- Finite high-qu natural text & synthetic cal factors and lannens

## Data & Protocol
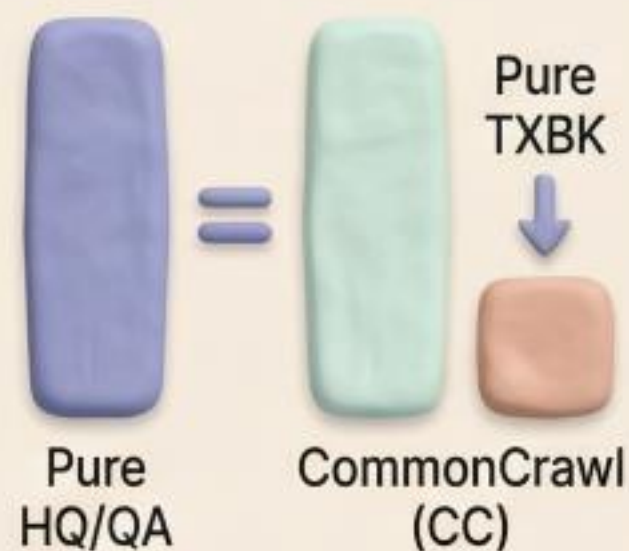
Up to 200B Tokens

Standardized Llama 3 Architecture

- Up to ~200B Tokens
- Standardized Llama 3 Architecture

## Training Mixtures

100% Natural

100% Synthetic

Mixed (67%/33%, 33%/67%)

- Syntheticic empirical study
- Finite health quality contents benefits ehectetic concerns

# Irreducible Loss & Model Collapse Evidence

## Irreducible Loss Estimations ($E$)

Synthetic Mixtures (except pure QA)

Lower Irreducible Loss than CommonCrawl only

33% HQ + 67% CC

Lowest Projected Irreducible Loss

## Model Collapse Evidence

Rephrased Data

No Degradation

Textbook-Style Data

Shows Patterns of Model Collapse

# Low-Level Statistical Analysis

## Unigram Frequency Insights

Training-test mismatch from rare unigrams (e.g., ' ', 'hvor') in training sets causes higher loss.

No single training set offers complete coverage.

Synthetic data slightly shrinks unigram distribution vs. CommonCrawl.

## Key Takeaways

CommonCrawl has widest coverage & lowest KL-divergence, but doesn't yield superior performance.

Good training mixtures depend on complex diversity-quality trade-offs beyond simple similarity.

Performance depends on complex diversity-quality trade-offs, not just coverage.

# Key Research Findings & Practical Guidance

## Strategic Mixture Accelerates 5-10x

Strategic incorporation of specific synthetic data types can accelerate pre-training convergence 5-10x.

Convergence

1/3 Rephrased
2/3 Natural

At Larger Budgets

## Synthetic Type Matters

Impact is highly dependent on synthetic data type; rephrased alone not faster than natural, textbook alone results in higher

Rephrased Alone = Not Faster

Textbook Alone = Higher Loss

## Optimal Ratio (~30%)

~30%

Rephrased Data

"Good" ratios vary with type, model scale, and budget, converging to ~30% for rephrased data.

Type    Model Scale    Budget

## Generator Scale ≠ Better Data

Larger Generator                    ~8B Model

Larger generator models don't necessarily yield better pre-training data than ~8B models.

# Limitations & Future Directions

Demystifying Synthetic Data in LLM Pre-training

## Study Limitations

This study has several limitations in scope, evaluation, and scale that must be considered.

Scope
Evaluation
Scale

## Limited Data Types

Limited scope to three synthetic data types:
- HQ (High-Quality)

HQ · QA · TXBK

- TXBK (Textbook)

Other types not explored.

## Reliance on Perplexity/Loss

- Heavy reliance on perplexity/loss metrics.
- Lacks in-depth human evaluation for qualitative assessment.
- Metrics may not capture true capabilities.

Perplexity

Human Evaluation?

## Single Pre-training Stage

- Focuses on a single pre-training stage.
- No analysis of long-term or multi-generational effects.
- Potential for degradation over time not studied.

Long-Term?

Future Generations

## Scale Constraints

- Models scaled up to 3B parameters.
- Tokens limited to 200B.
- Requires validation at frontier scales (e.g., >100B parameters, trillions of tokens).

Frontier Scale >100B

3B

## Impact of Tokenizers

- Impact of tokenizers on vocabulary coverage not fully explored.
- Suboptimal tokenization can limit model performance.
- Different tokenizers may yield different results.

Tokenizer

## Future Directions

Proposing future research directions to address current limitations and advance the field.

## Key Areas for Future Research

Targeted Generation
Develop more targeted synthetic data generation techniques.

Dynamic Mixing
Explore dynamic mixing strategies.

Long-Term Impact
Rigorous evaluation of long-term impacts on diverse capabilities.

Generator Characteristics
Identify key generator characteristics beyond just size.

# Ethical Considerations & Conclusion

## Ethical Risks & Mitigation

**Bias Propagation Risks**
(from Generator Models)

**Factual Accuracy Concerns**
(Misinformation)

**Reduced Diversity**
(Over-reliance)

**Importance of Transparency**
(Mitigation)

## A Nuanced Trade-off & Conclusion

Trade-off — Empirically-Informed

**Nuanced Trade-off**
(Requires Careful Deployment)

**Empirical Deployment**
(Not a Universal Solution)

**Substantial Benefits**
(Strategic Mixing Mixing Despite Theoretical Concerns)