

# Simplicial Sets: Designing the Geometric Transformer

Sridhar Mahadevan, Adobe Research and U.Mass, Amherst

# Recap: Lectures 1-6

- **Categories**: collection of **objects** and **arrows**
- **Functors**: structured mapping between categories
- **Natural transformations**: structured mappings between functors
- **Yoneda Lemma**: Representable set-valued functors
- **Limits and colimits**: Turning diagrams into computational objectives
- **Functor categories**: Each object is a functor

# Diagrammatic Backpropagation: Deep learning with (co)limits

- ✦ The textbook *Categories for AGI* has several chapters devoted to this week's topics
  - ✦ Chapters 5-10 (~100 pages!)
  - ✦ It is the longest section in the book!
- ✦ We will study these concepts further in the course
  - ✦ This week's lectures are a high-level summary

# GAIA: Generative AI Architecture

## GAIA: CATEGORICAL FOUNDATIONS OF GENERATIVE AI\*

A PREPRINT

Sridhar Mahadevan

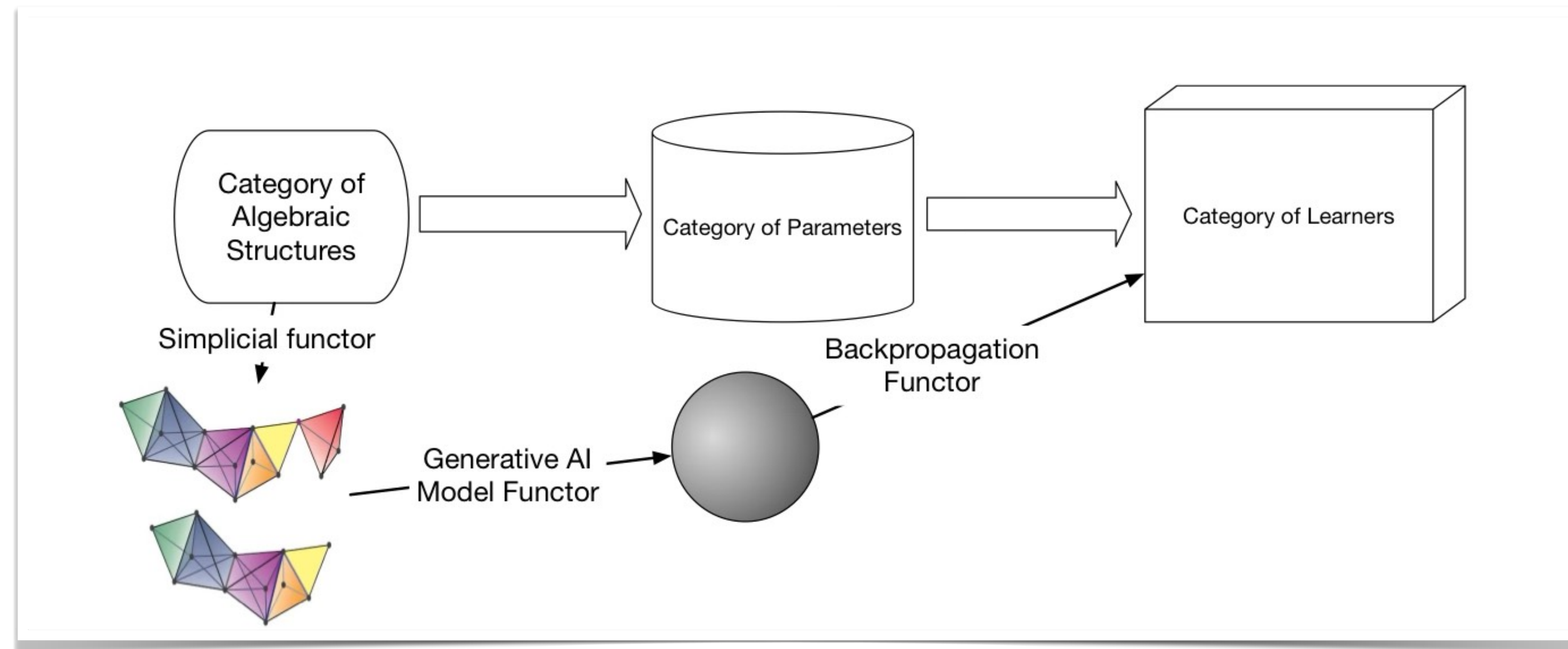
Adobe Research and University of Massachusetts, Amherst  
smahadev@adobe.com, mahadeva@umass.edu

March 1, 2024

### ABSTRACT

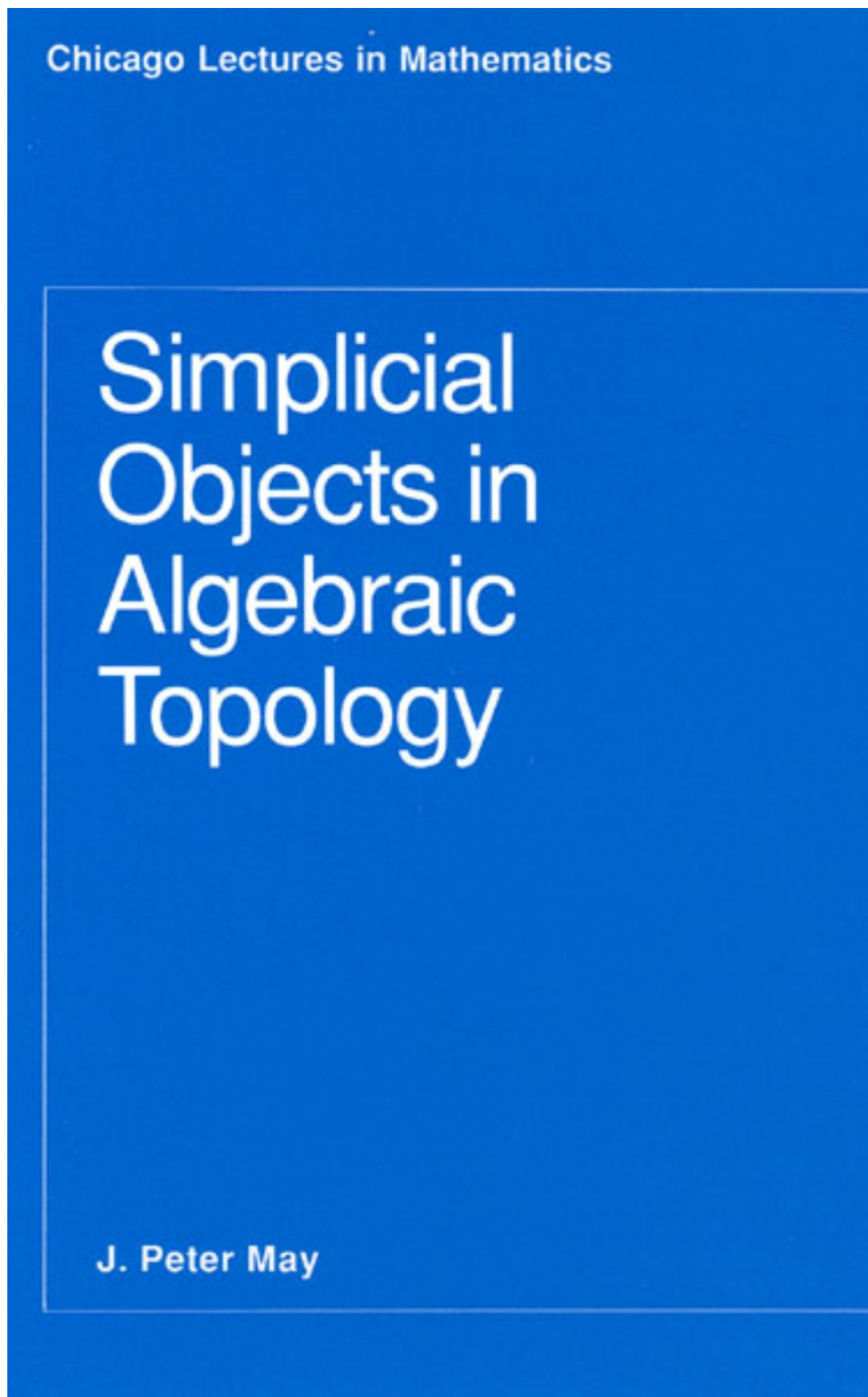
In this paper, we explore the categorical foundations of generative AI. Specifically, we investigate a Generative AI Architecture (GAIA) that lies beyond backpropagation, the longstanding algorithmic workhorse of deep learning. Backpropagation is at its core a compositional framework for (un)supervised learning: it can be conceptualized as a sequence of modules, where each module updates its parameters based on information it receives from downstream modules, and in turn, transmits information back to upstream modules to guide their updates. GAIA is based on a fundamentally different *hierarchical model*. Modules in GAIA are organized into a simplicial complex. Each  $n$ -simplicial complex acts like a manager of a business unit: it receives updates from its superiors and transmits information back to its  $n + 1$  subsimplicial complexes that are its subordinates. To ensure this simplicial generative AI organization behaves coherently, GAIA builds on the mathematics of the higher-order category theory of simplicial sets and objects. Computations in GAIA, from query answering to foundation model building, are posed in terms of lifting diagrams over simplicial objects. The problem of machine learning in GAIA is modeled as “horn” extensions of simplicial sets: each sub-simplicial complex tries to update its parameters in such a way that a lifting diagram is solved. Traditional approaches used in generative AI using backpropagation can be used to solve “inner” horn extension problems, but addressing “outer horn” extensions requires a more elaborate framework.

At the top level, GAIA uses the simplicial category of ordinal numbers with objects defined as  $[n], n \geq 0$  and arrows defined as weakly order-preserving mappings  $f : [n] \rightarrow [m]$ , where  $f(i) \leq f(j), i \leq j$ . This top-level structure can be viewed as a combinatorial “factory” for constructing, manipulating, and destructing complex objects that can be built out of modular components defined over categories. The second layer of GAIA defines the building blocks of generative AI models as universal coalgebras over categories that can be defined using current generative AI approaches, including Transformers that define a category of permutation-equivariant functions on vector spaces, structured state-space models that define a category over linear dynamical systems, or image diffusion models that define a probabilistic coalgebra over ordinary differential equations. The third layer in GAIA is a category of elements over a (relational) database that defines the data over which foundation models are built. GAIA formulates the machine learning problem of building foundation models as extending functors over categories, rather than interpolating functions on sets or spaces, which yields canonical solutions called left and right Kan extensions. GAIA uses the metric Yoneda Lemma to construct universal representers of objects in non-symmetric generalized metric spaces. GAIA uses a categorical integral calculus of (co)ends to define two families of generative AI systems. GAIA models based on coends correspond to topological generative AI systems, whereas GAIA systems based on ends correspond to probabilistic generative AI systems.

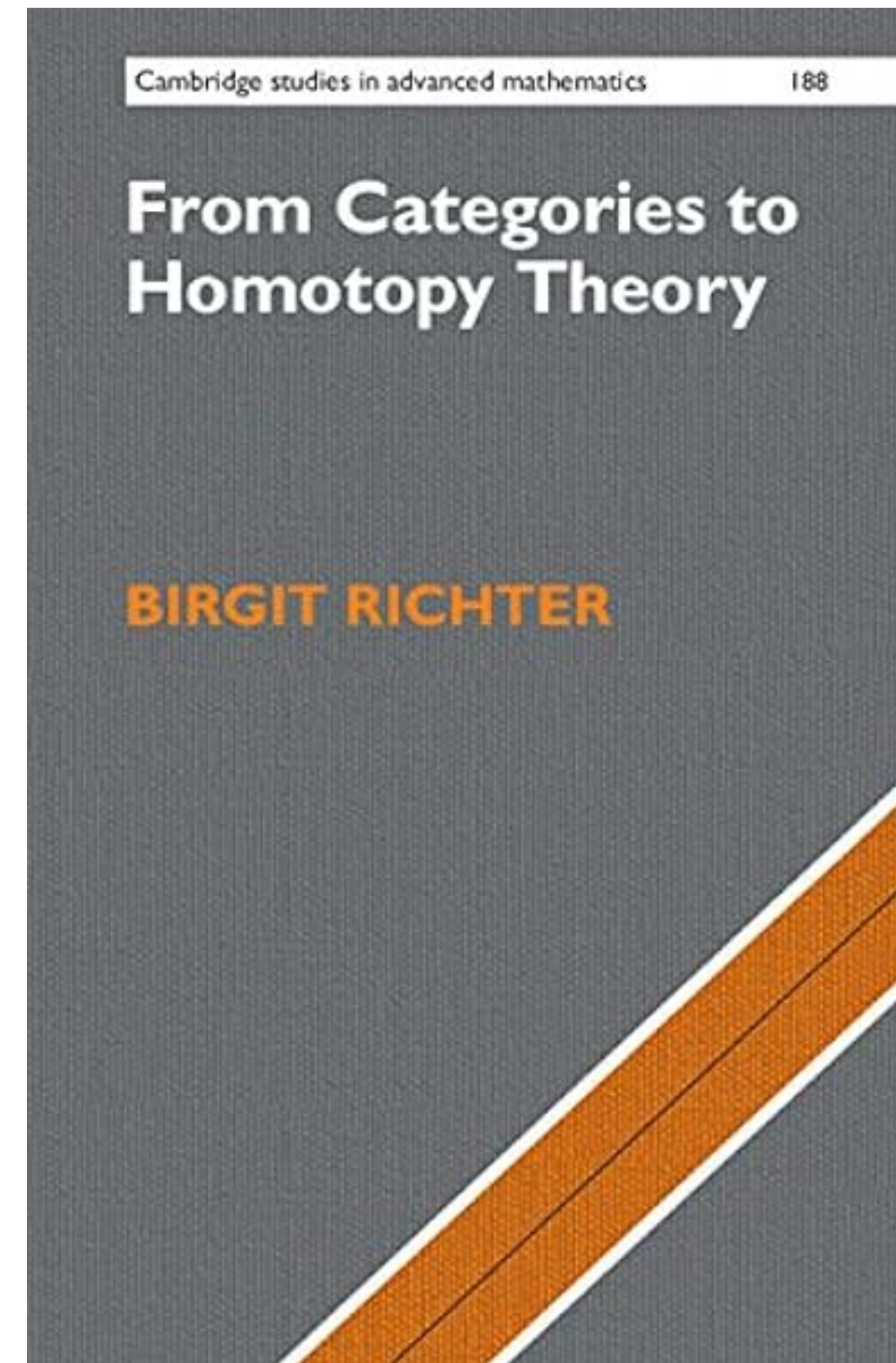
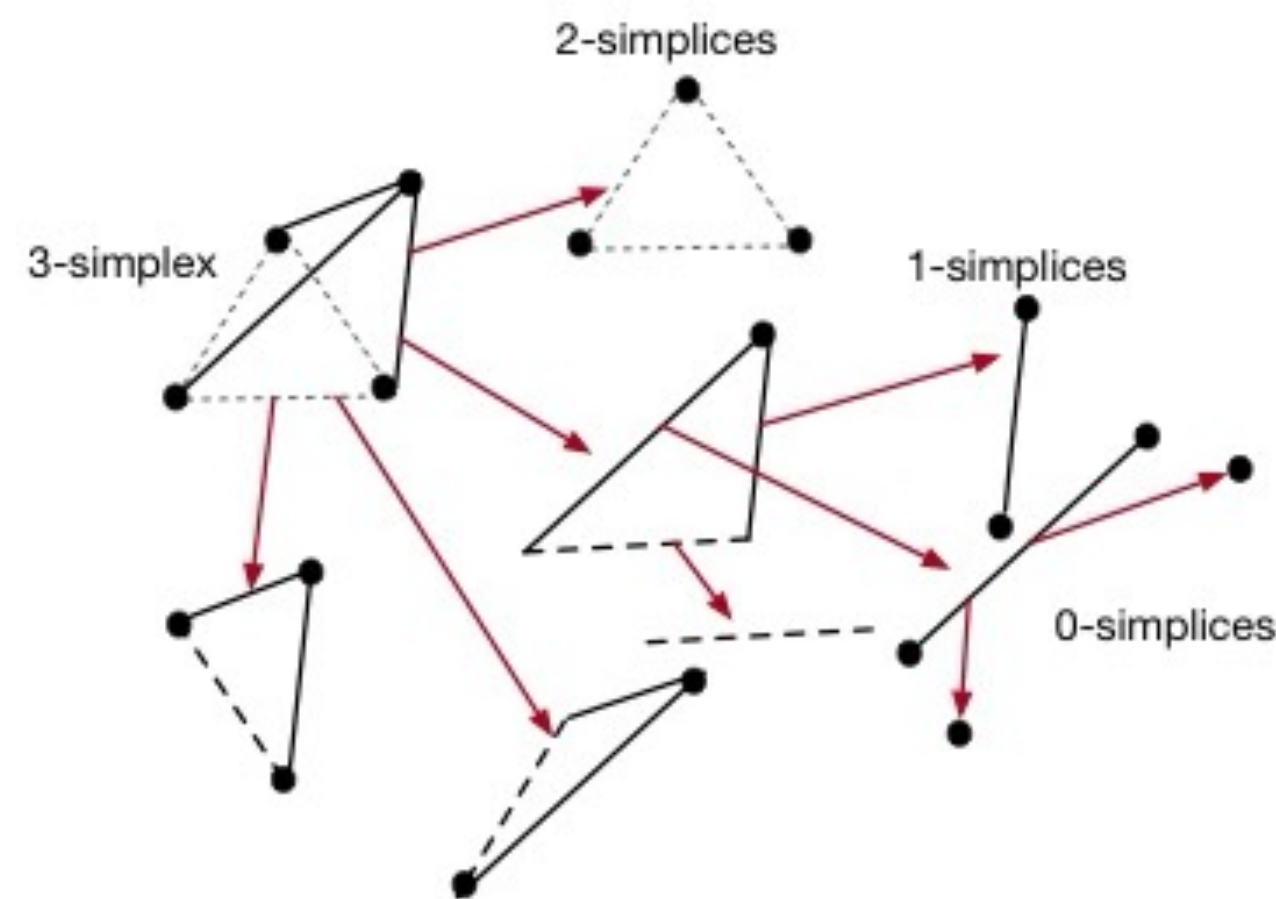
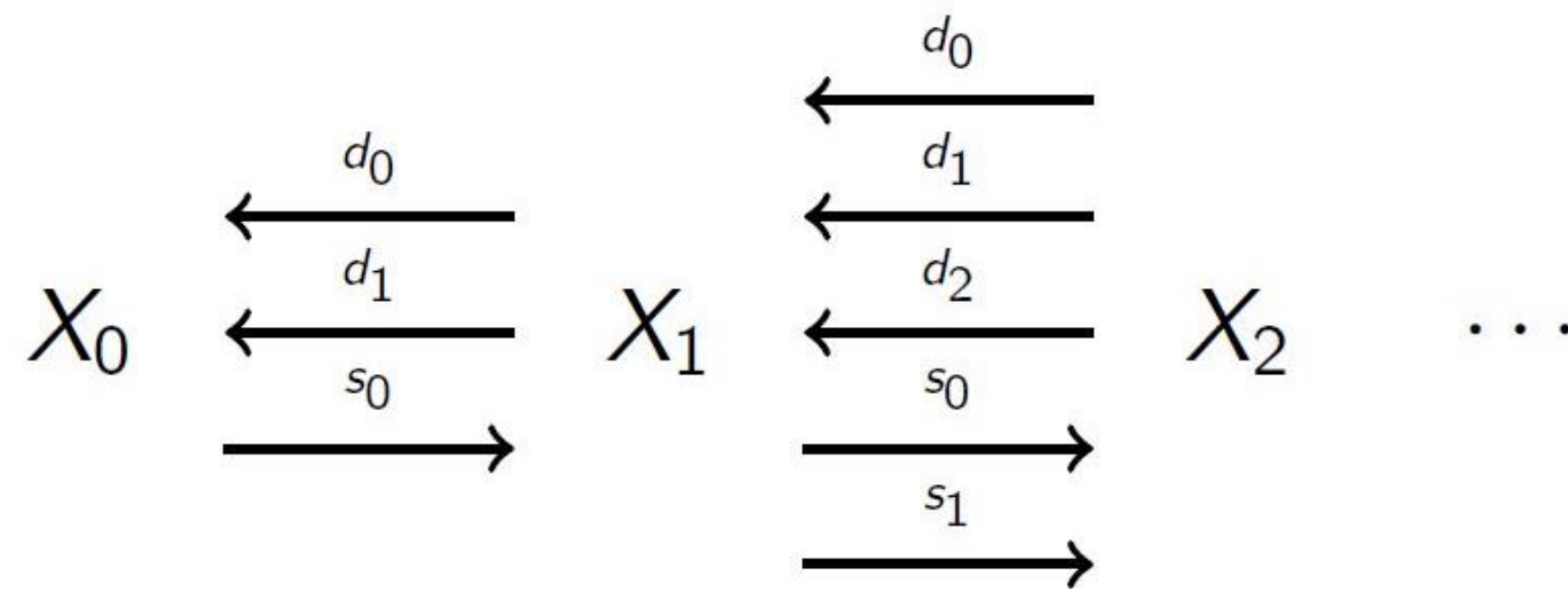


**Keywords** Generative AI · Foundation Models · Higher-Order Category Theory · Machine Learning

# Simplicial Sets and Objects



Combinatorial model of spaces



# Simplicial Category $\Delta$

- **Objects:** ordinal numbers

- $[n] = \{0, 1, \dots, n - 1\}$

- **Arrows:**

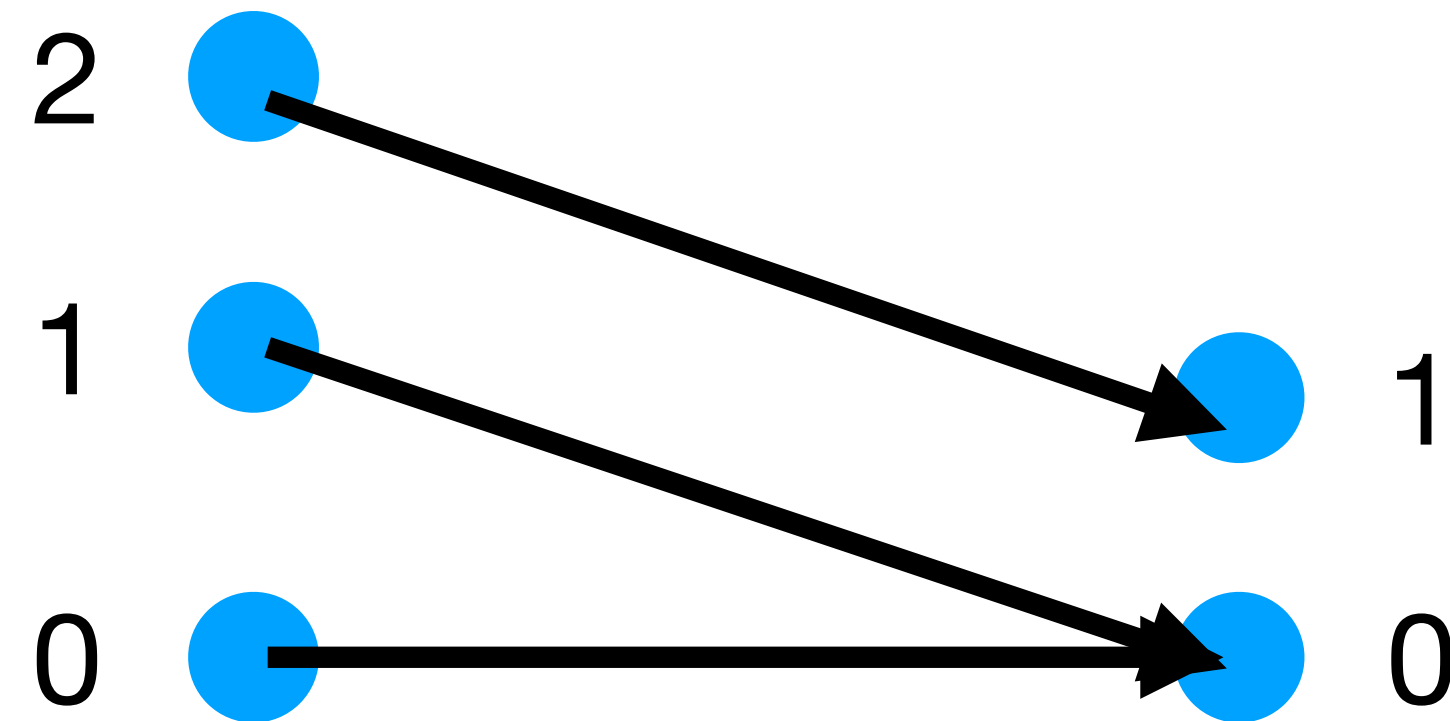
- $f : [m] \rightarrow [n]$

- If  $i \leq j$ , then  $f(i) \leq f(j)$

- All morphisms can be built out of primitive injections/surjections

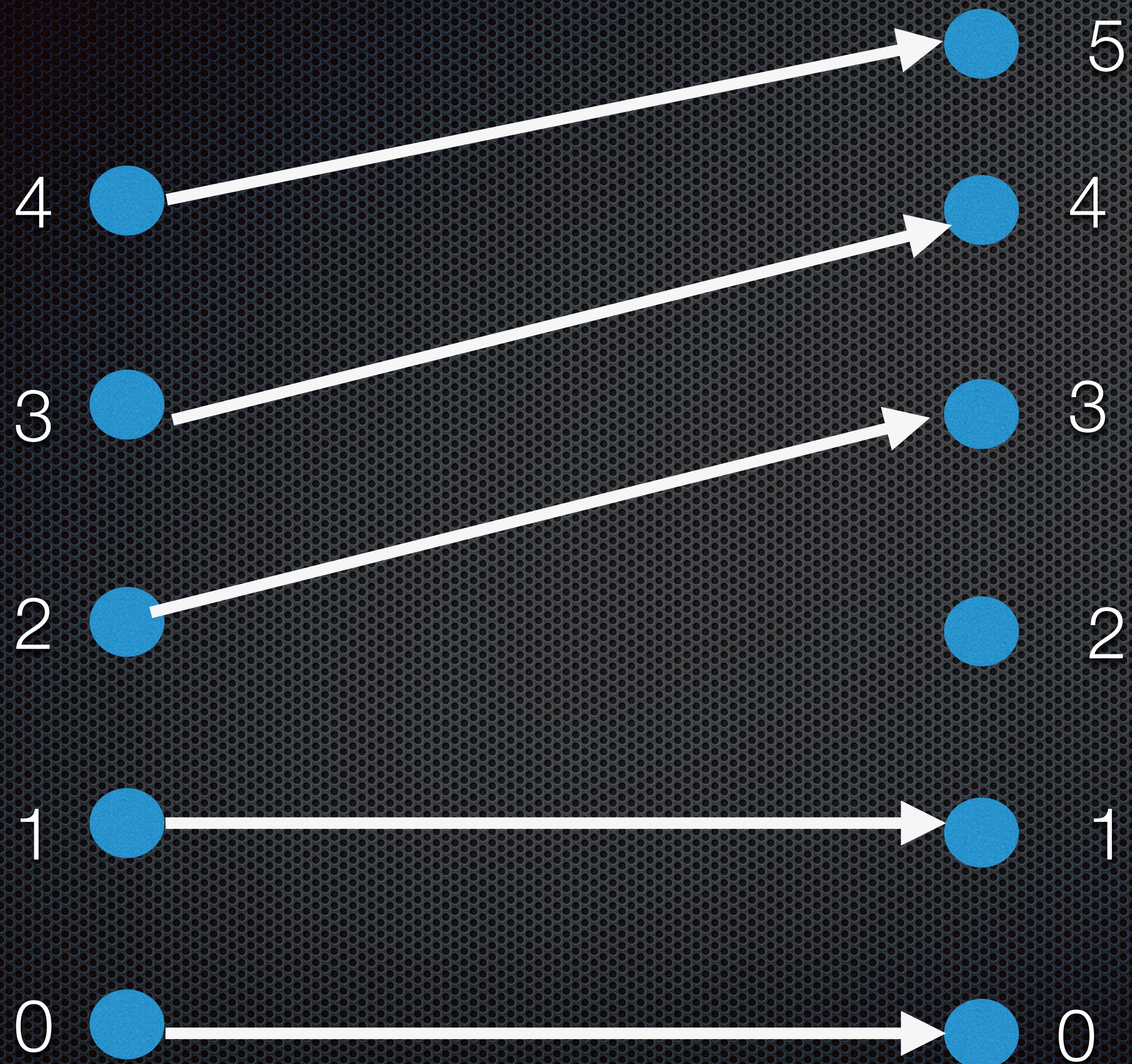
- $\delta_i : [n] \rightarrow [n + 1]$  : injection skipping  $i$

- $\sigma_i : [n] \rightarrow [n - 1]$ , surjection repeating  $i$

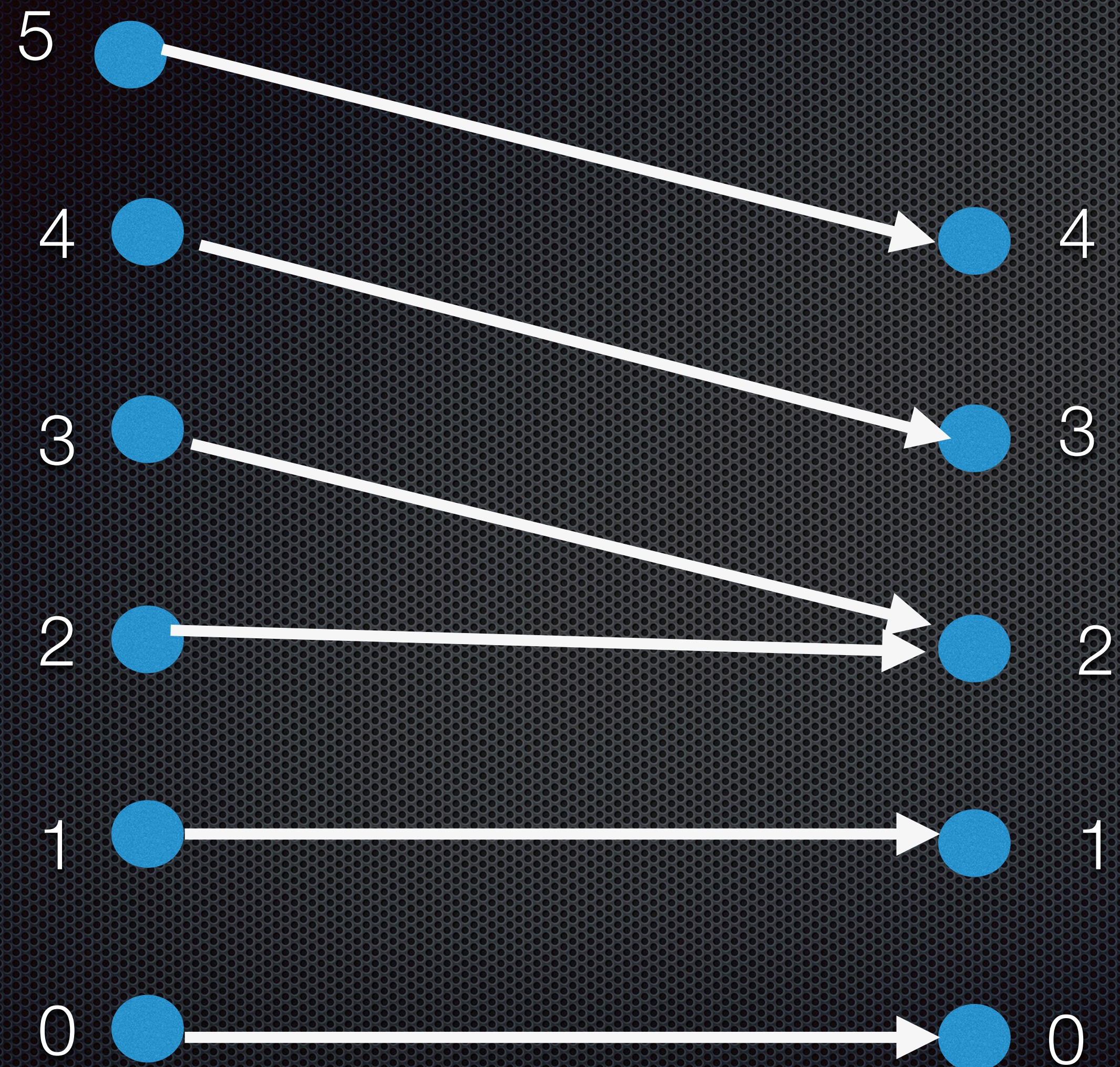


# Primitive Injections

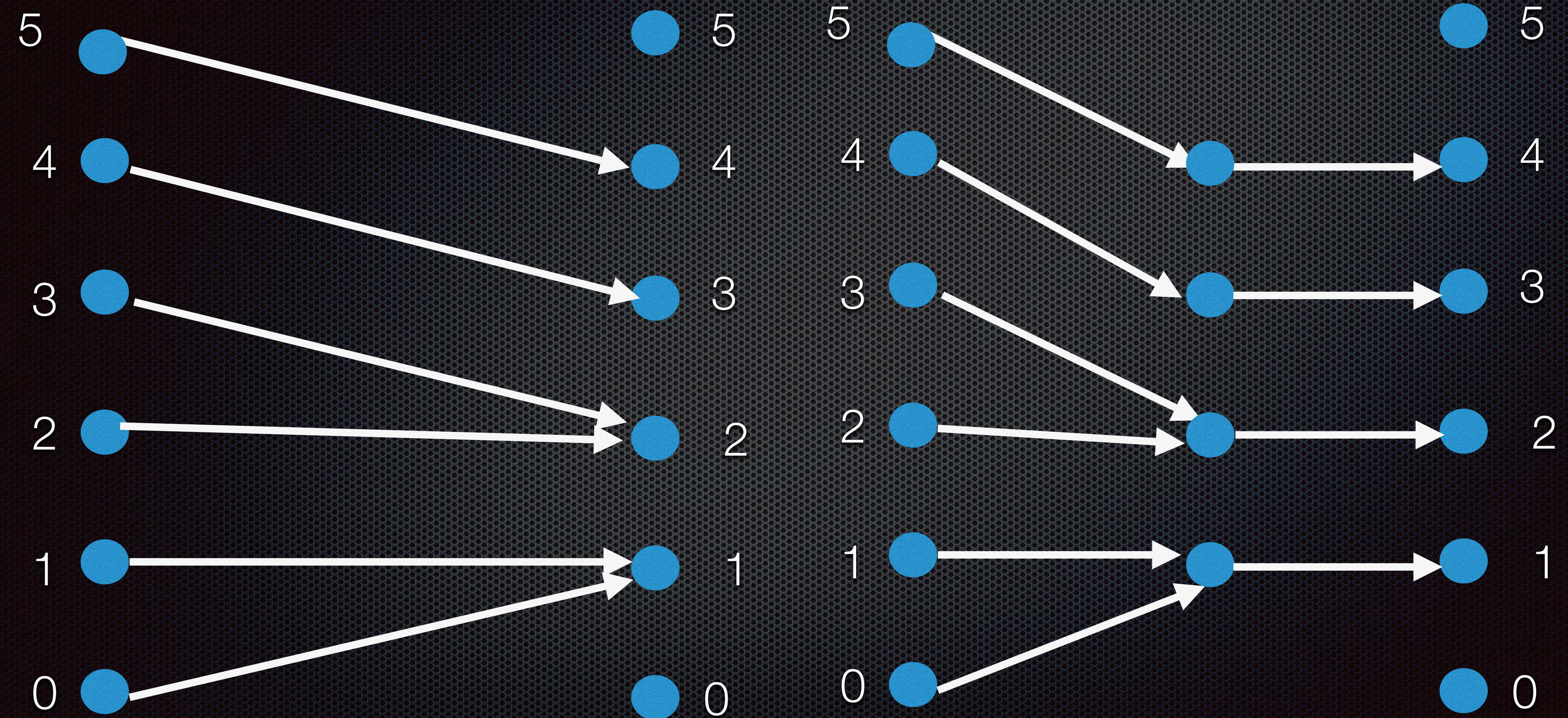
$$\delta_2 : [n] \rightarrow [n + 1]$$



# Primitive Surjections: $\sigma_2 : [n + 1] \rightarrow [n]$



# Every arrow is composition of primitive surjections + injections



# Yoneda lemma for simplicial sets

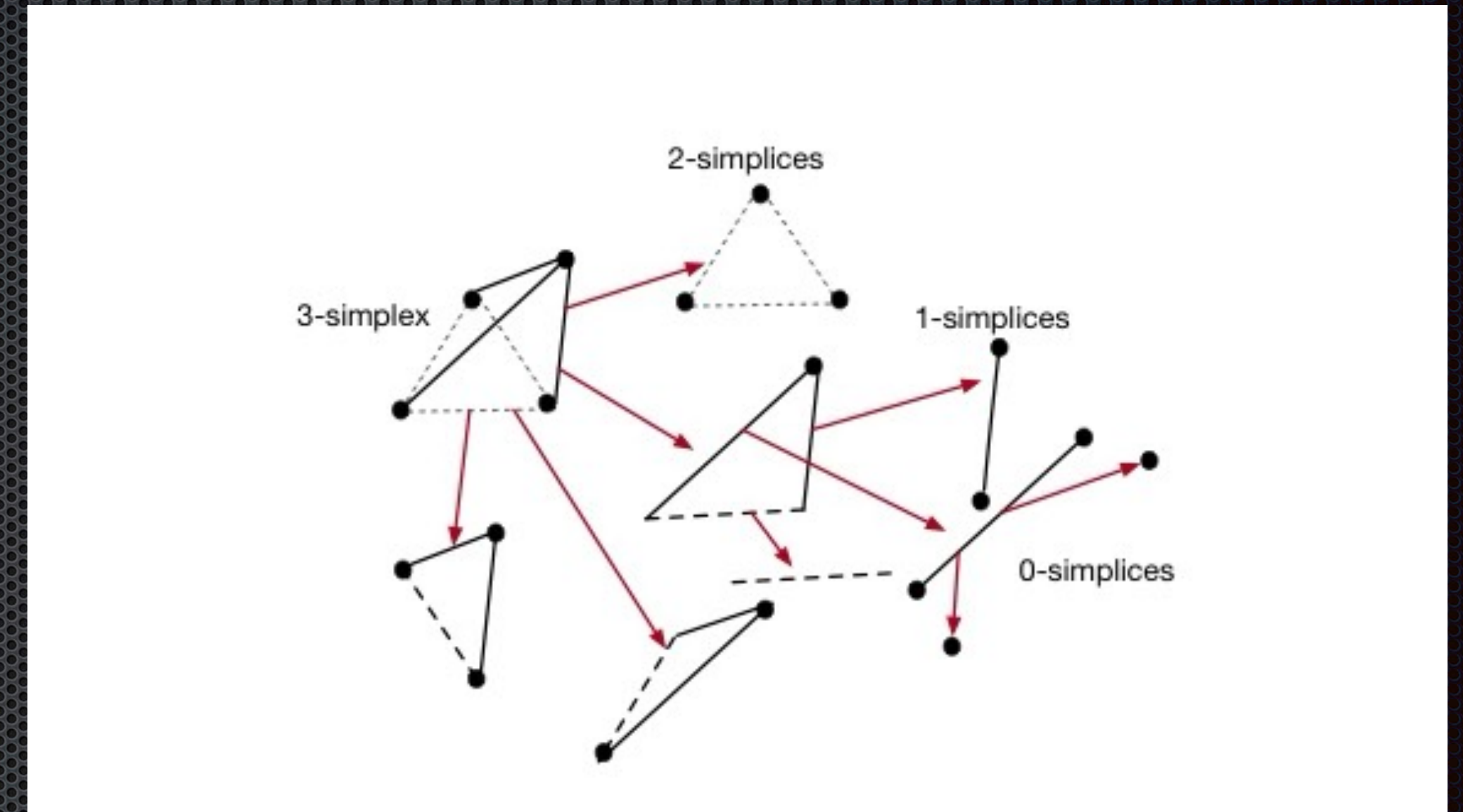
$$\mathbf{sSet}: \Delta \rightarrow \mathbf{Set}^{\Delta^{op}}$$

$$\Delta^n = \Delta(-, [n])$$

$$\Delta_k^n = \Delta([k], [n])$$

$$d_i : f : [k] \rightarrow [n] \mapsto [k-1] \xrightarrow{\delta_i} [k] \xrightarrow{f} [n]$$

$$s_i : f : [k] \rightarrow [n] \mapsto [k+1] \xrightarrow{\sigma_i} [k] \xrightarrow{f} [n]$$

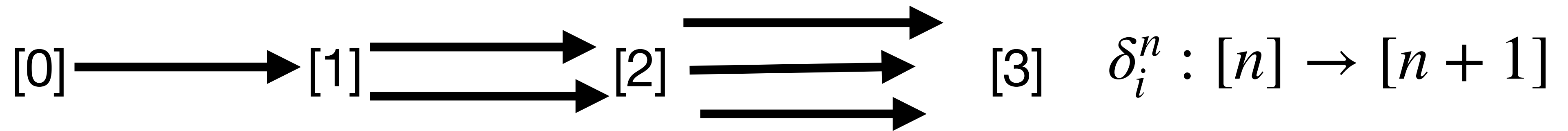


# Nerve of a Category

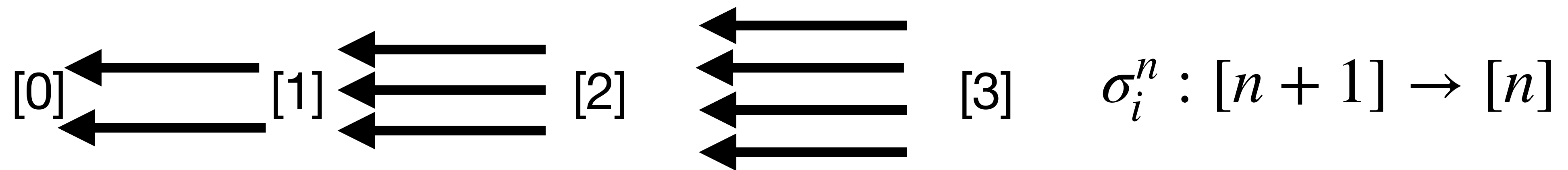
- A simplicial set is a contravariant functor mapping the simplicial category to the category of sets
- Any category can be mapped onto a simplicial set by constructing its nerve

$$M_n(\mathcal{C}) = \{c_0 \xrightarrow{f_1} c_1 \xrightarrow{f_2} c_1 \cdots \xrightarrow{f_n} c_n \mid c_i \text{ is an object and } f_i \text{ is an arrow}\}$$

# Simplicial Sets: Contravariant Functors



$$\delta_i^n : [n] \rightarrow [n + 1]$$



$$\sigma_i^n : [n + 1] \rightarrow [n]$$

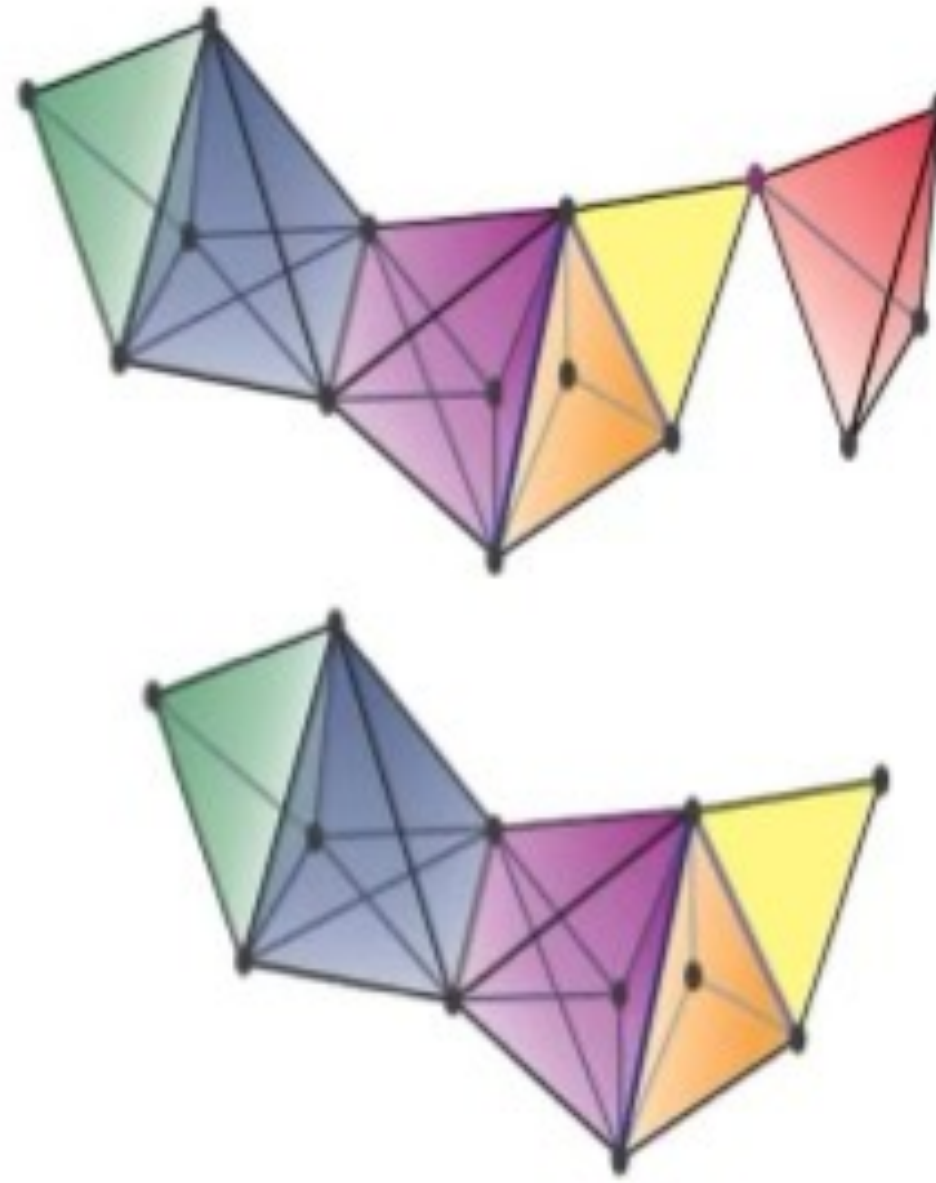
$$X_n : [n] \rightarrow X : \Delta^{op} \rightarrow X$$

# Simplicial Sets vs. Categories

- Any category can be embedded faithfully into a simplicial set using its nerve
- The embedding is full and faithful (perfect reconstruction)
- Unfortunately, the converse is not possible
- Given a simplicial set, the left adjoint functor that maps it into a category is lossy!

### Simplicial framework for generative AI

Simplicial learning is based on extension problems of inner and outer "horns" of simplicial objects



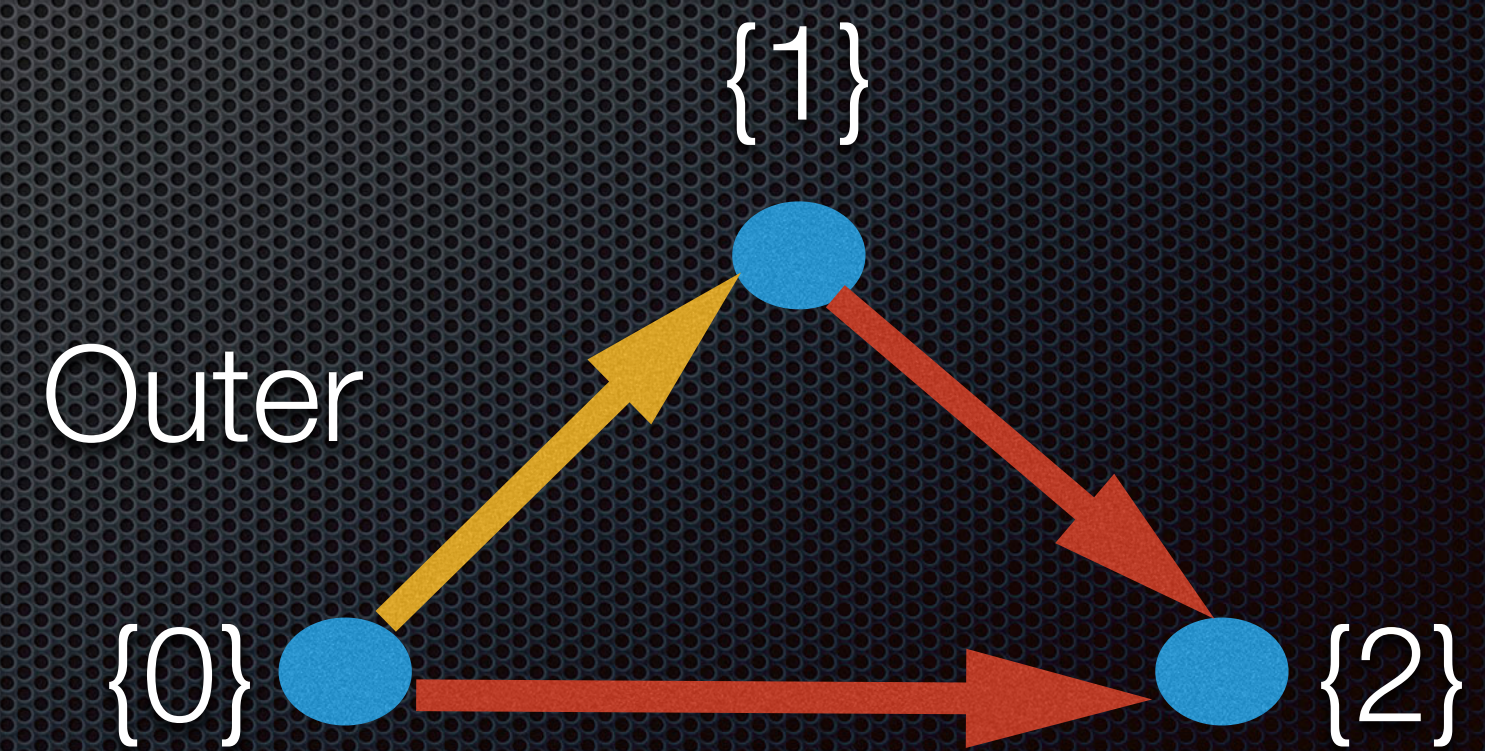
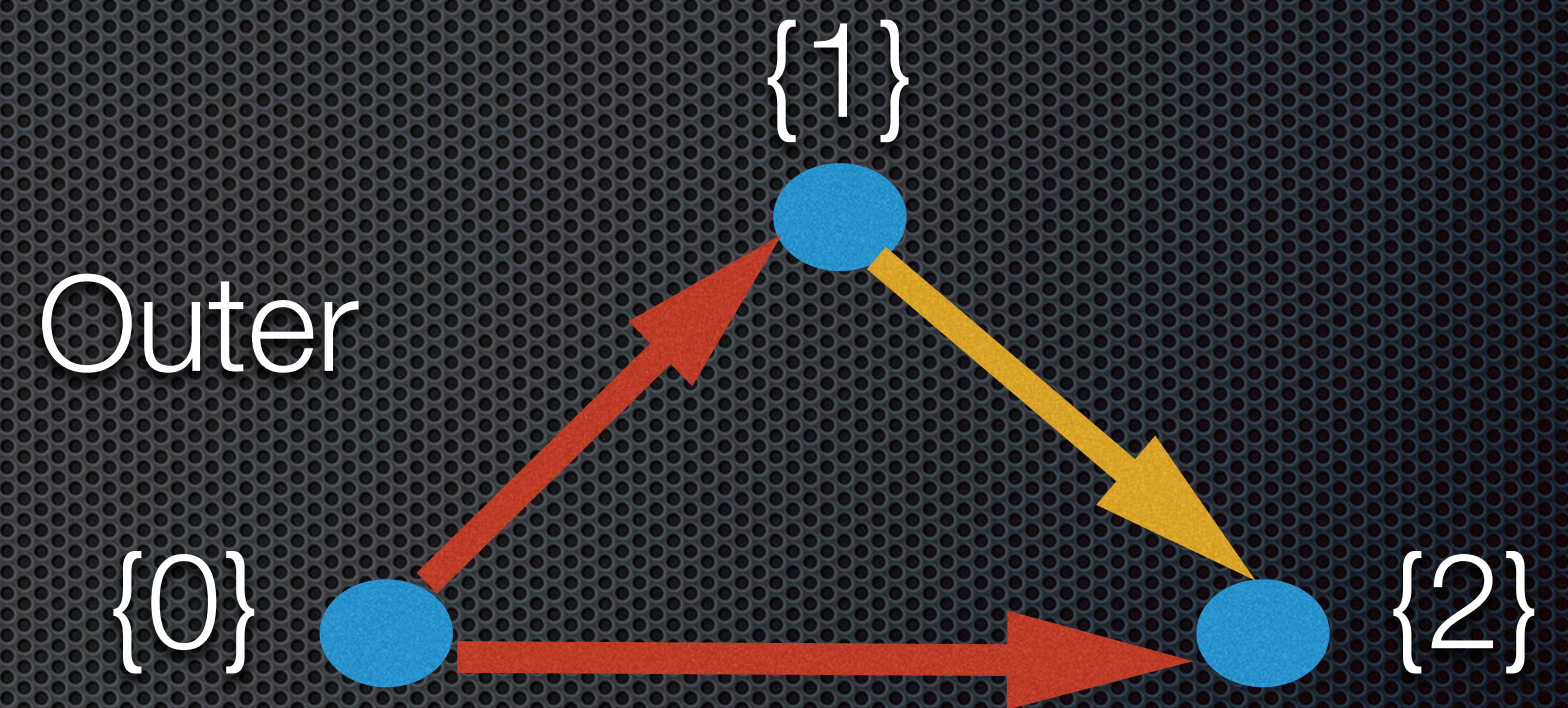
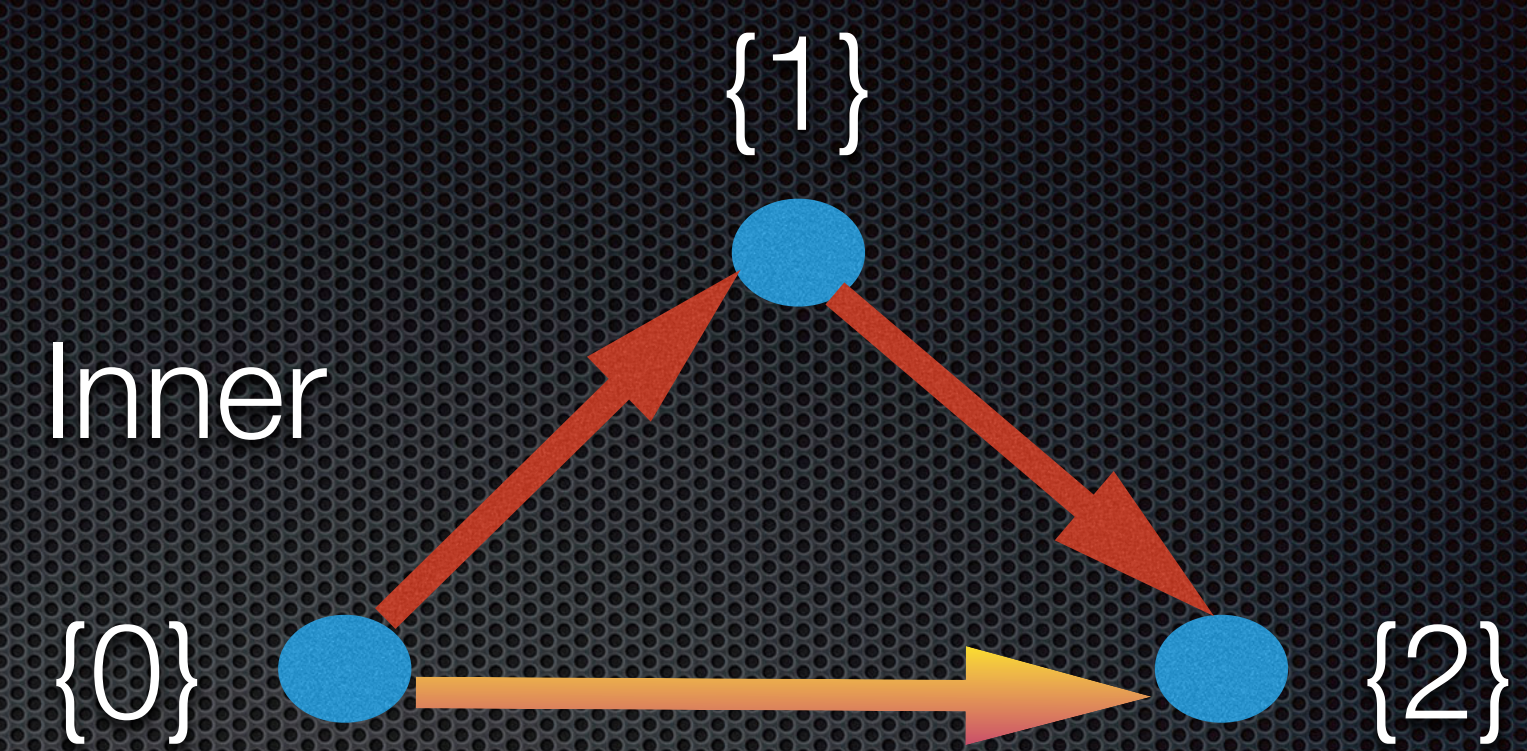
Each directed edge defines a morphism that represents a generative AI method

Each collection of simplices can be "glued" on to compatible simplices through "ports" that define the components of the simplex.

# Horns of Simplicial Sets

A “horn” is a subset of a simplicial set

The “extension” problem is to “fill in” missing faces

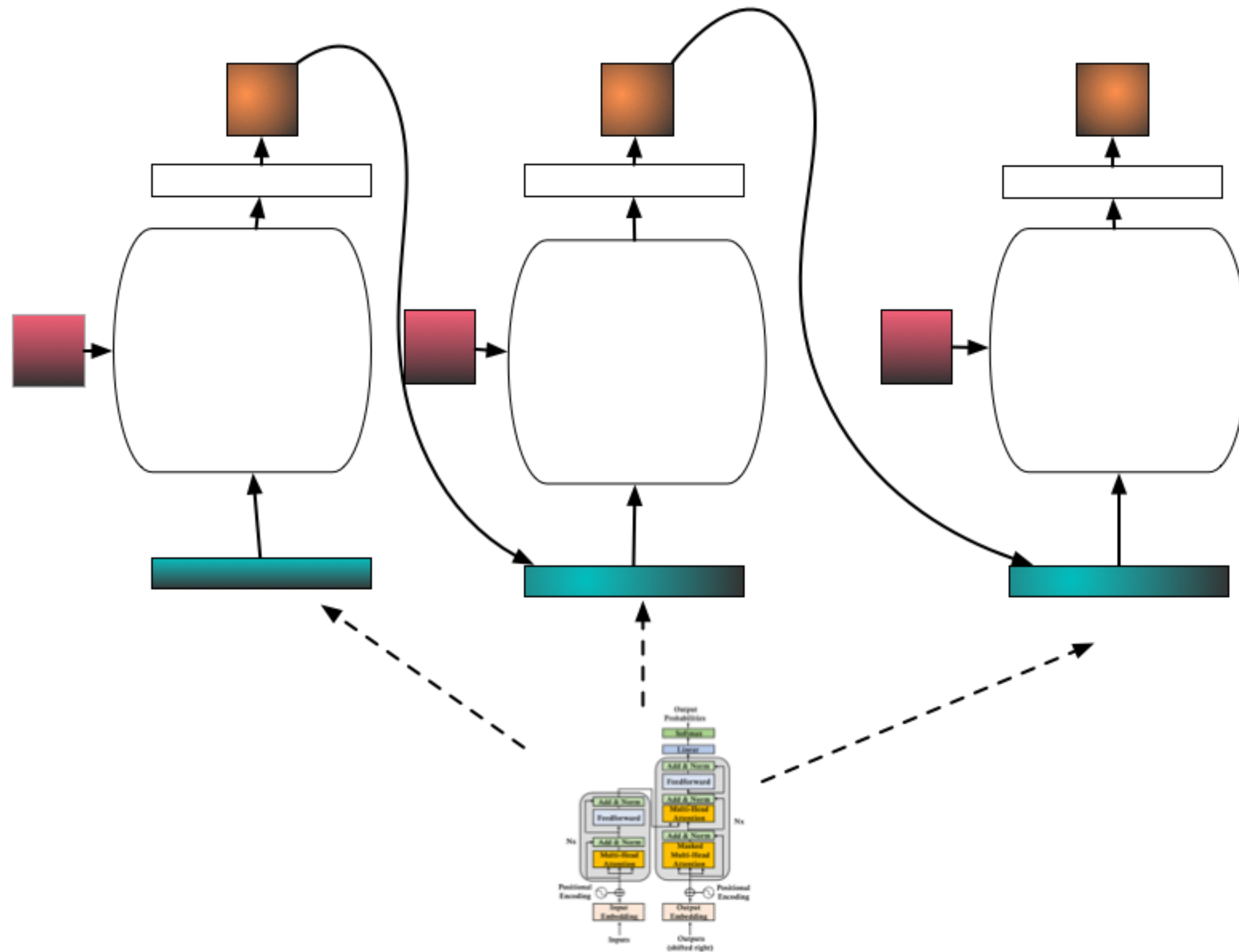


# Lifting Diagrams for Simplicial Sets

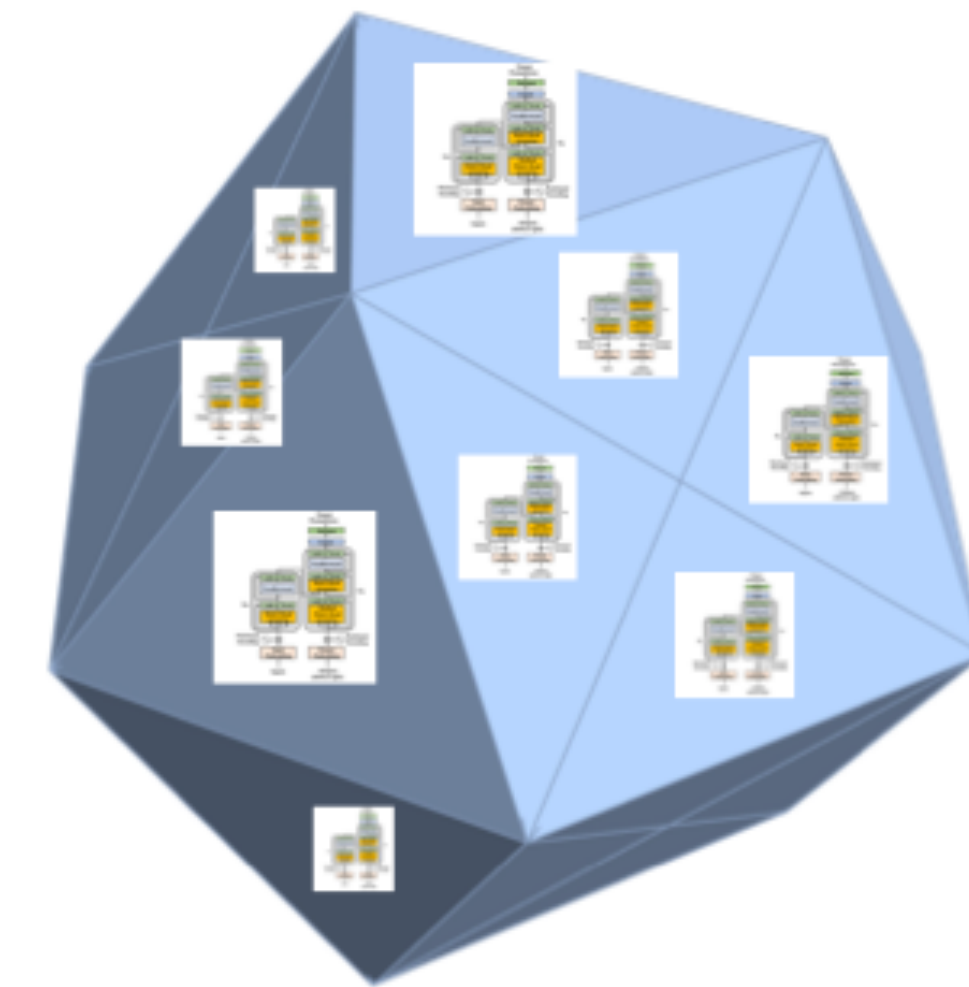
Recall that a lifting diagram  
was used as a way to model  
inference in databases.

# Simplicial Generative AI

Sequential Generative AI Model



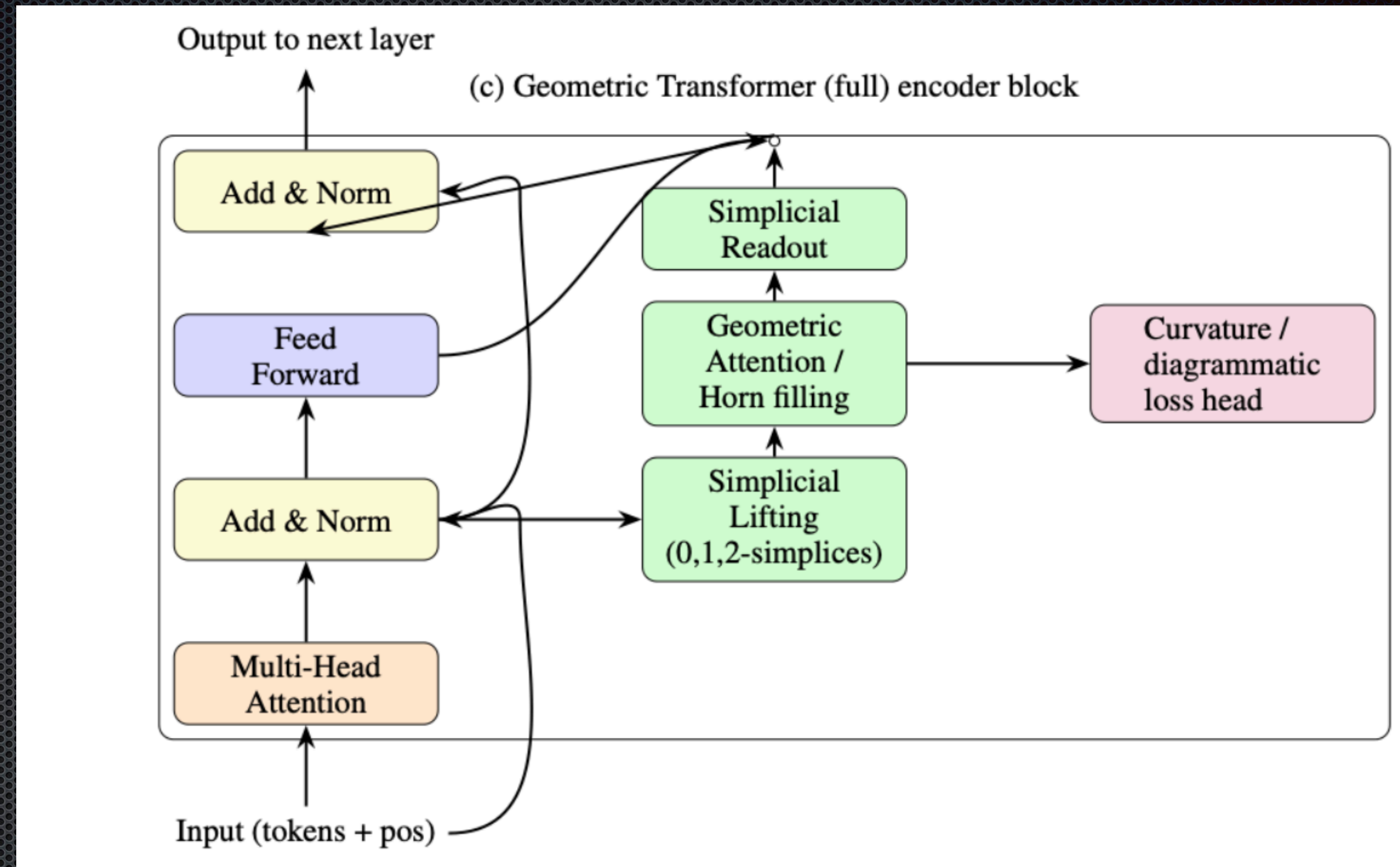
Simplicial Complex of Gen AI Models in GAIA



# Geometric Transformers

GTs introduce explicit geometric inductive biases

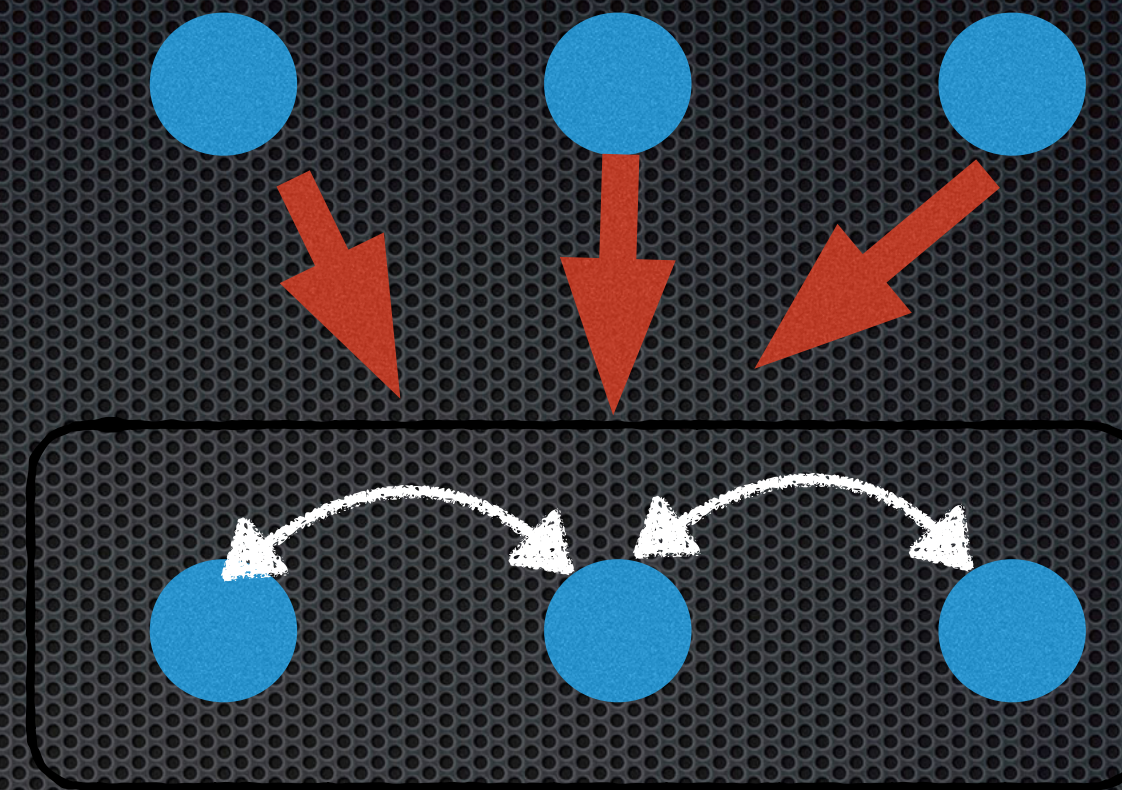
Token embedding is geometrically transported into a simplicial space



[Mahadevan, Categories for AGI, 2026]

# GT-Lite

Uses a simple 1D convolution layer to “smooth” out token embeddings on adjacent tokens



1D convolution is the colimit  
in a vector space embedding

# Simplicial Message Passing: GT-Lite

Uses a discrete diffusion  
process

Attenuates local  
inconsistencies

---

Algorithm 9: GT-Lite Local Geometric Smoothing Operator  $C$  (1D Conv + Residual + LN)

**Input:** Token states  $z \in \mathbb{R}^{B \times L \times d}$ ; convolution weights  $\theta_C$ ; smoothing scale  $\alpha$  (e.g., 0.2); LayerNorm parameters  $\theta_{LN}$

**Output:** Smoothed token states  $C(z) \in \mathbb{R}^{B \times L \times d}$

- 1: // 1) Convert to Conv1d layout and apply local neighborhood mixing
  - 2:  $z^\top \leftarrow \text{transpose}(z, (B, L, d) \rightarrow (B, d, L))$
  - 3:  $u^\top \leftarrow \text{Conv1D}_{\theta_C}(z^\top) \triangleright$  kernel size 3, padding 1, output shape  $(B, d, L)$
  - 4:  $u \leftarrow \text{transpose}(u^\top, (B, d, L) \rightarrow (B, L, d))$
  
  - 5: // 2) Residual update and normalization
  - 6:  $z' \leftarrow z + \alpha \cdot u$
  - 7:  $C(z) \leftarrow \text{LayerNorm}_{\theta_{LN}}(z')$
  - 8: **return**  $C(z)$
-

# Diagrammatic Backpropagation

- ❖ Central idea: penalize deviations from compositionality
- ❖ Propagate gradients around loops!
- ❖ By explicitly measuring the lack of compositionality, we get a novel error signal!

[Mahadevan, Categories for AGI, 2026]

Algorithm 7: Diagrammatic Backpropagation Training Loop (code-faithful)

**Input:** Model  $M_\theta$  (Transformer or GT variant); dataset iterator  $\mathcal{B}$  yielding minibatches  $(s, y)$ ; optimizer ADAMW; steps  $T$ ; warmup steps  $T_w$ ; weights  $\lambda_{\text{db}}, \lambda_{\text{geom}}$ ; clipping norm  $c$ ; small constant  $\epsilon$

**Output:** Trained parameters  $\theta$

```
1: Initialize: set global step  $t \leftarrow 0$ 
2: for  $t = 1$  to  $T$  do
3:   Sample minibatch  $(s, y) \sim \mathcal{B} \triangleright s$ : tokens (or inputs),  $y$ : labels/targets

4:   // 1) Forward pass for the primary task
5:    $\hat{y} \leftarrow M_\theta(s)$ 
6:    $\mathcal{L}_{\text{task}} \leftarrow \text{LossTask}(\hat{y}, y) \triangleright$  e.g., cross-entropy for LM

7:   // 2) Optional DB control loss (commutator-energy terms)
8:    $\mathcal{L}_{\text{db}} \leftarrow 0$ 
9:   if  $\lambda_{\text{db}} > 0$  then
10:     $(\text{Obs}, \{\text{Obs}_\ell\}) \leftarrow \text{COMMProbe}(M_\theta, s; \epsilon) \triangleright$  Algorithm 6
11:     $\mathcal{L}_{\text{db}} \leftarrow \text{Obs} \triangleright$  scalar average over layers; can also weight
    layers/heads
12:   end if

13:   // 3) Optional diagrammatic curvature loss (triangle/square residuals)
14:    $\mathcal{L}_{\text{geom}} \leftarrow 0$ 
15:   if  $\lambda_{\text{geom}} > 0$  and  $(s, y)$  includes a compiled diagram/simplicial instance then
16:     $\triangleright$  For diagram benchmarks,  $s$  contains  $(\mathcal{D}, x)$  with constraints.
17:     $\mathcal{L}_{\text{geom}} \leftarrow \text{CURVATUREENERGY}(M_\theta, s) \triangleright$  triangle/square
    residual energy; cf. v1 Alg. (learned energy)
18:   end if

19:   // 4) Warmup schedule for auxiliary terms (optional but common)
20:    $\lambda_t^{\text{db}} \leftarrow \lambda_{\text{db}} \cdot \min(1, t/T_w)$ 
21:    $\lambda_t^{\text{geom}} \leftarrow \lambda_{\text{geom}} \cdot \min(1, t/T_w)$ 

22:   // 5) Total loss and parameter update
23:    $\mathcal{L} \leftarrow \mathcal{L}_{\text{task}} + \lambda_t^{\text{db}} \mathcal{L}_{\text{db}} + \lambda_t^{\text{geom}} \mathcal{L}_{\text{geom}}$ 
24:    $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L} \triangleright$  implemented via AdamW step
25:   Clip gradients:  $\|\nabla_\theta \mathcal{L}\| \leq c$ 
26:   Optional: normalize selected parameter groups (e.g., edge transforms
    / relation embeddings)
27: end for
28: return  $\theta$ 
```

# Simplicial Message Passing: GT-Full

Uses a discrete diffusion  
process

Attenuates local  
inconsistencies

---

Algorithm 8: GT-Full Geometric Transport Operator  $G$  (Simplicial/Graph Message Passing)

**Input:** Token states  $z \in \mathbb{R}^{B \times L \times d}$ ; fixed adjacency  $\text{edge\_index} \in \mathbb{N}^{2 \times E_{\max}}$  over positions  $0, \dots, L_{\max} - 1$ ; relation ids  $\text{rel\_ids} \in \mathbb{N}^{E_{\max}}$ ; domain ids  $\text{dom\_ids} \in \mathbb{N}^{E_{\max}}$ ; depth  $K$ ; learnable parameters  $\theta_G$

**Output:** Transported token states  $G(z) \in \mathbb{R}^{B \times L \times d}$

```
1: // 1) Flatten batch and restrict edges to current sequence length
2:  $h \leftarrow \text{reshape}(z, (B \cdot L) \times d)$ 
3:  $m \leftarrow \{e \in [1..E_{\max}] : \text{edge\_index}[0, e] < L\}$   $\triangleright$  mask endpoints within  $L$ 
4:  $E \leftarrow \text{edge\_index}[:, m]$ ;  $r \leftarrow \text{rel\_ids}[m]$ ;  $dID \leftarrow \text{dom\_ids}[m]$ 

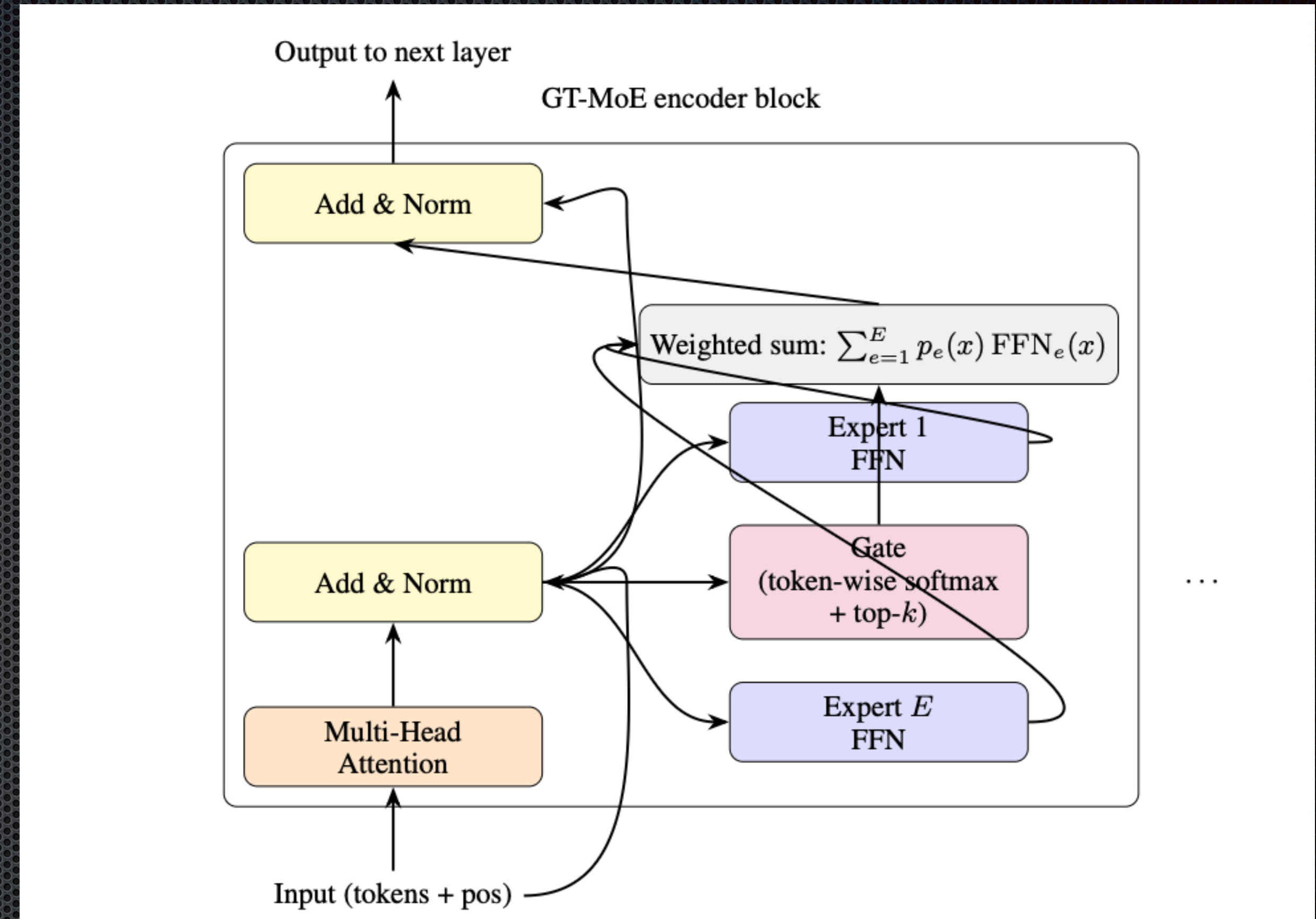
5: // 2) Iterate  $K$  rounds of geometric message passing
6: for  $k = 1$  to  $K$  do
7:   Initialize message accumulator  $M \leftarrow 0 \in \mathbb{R}^{(BL) \times d}$ 
8:   for each directed edge  $(u \leftarrow v)$  in  $E$  do
9:      $e \leftarrow (u, v)$ 
10:    Edge embedding:  $e_{\text{emb}} \leftarrow \text{Emb}_{\text{rel}}(r_e) + \text{Emb}_{\text{dom}}(dID_e)$ 
11:    Compute message:  $m_{v \rightarrow u} \leftarrow \phi_{\theta}(h_v, h_u, e_{\text{emb}})$   $\triangleright$  e.g., MLP on  $(h_v, h_u, e_{\text{emb}})$  or attention-style score
12:     $M_u \leftarrow M_u + m_{v \rightarrow u}$ 
13:  end for
14:  Degree normalize:  $M \leftarrow \text{DegNorm}(M, E)$ 
15:  Update rule:  $h \leftarrow \text{Update}_{\theta}(h, M)$   $\triangleright$  e.g., GRUCell/residual MLP + LayerNorm
16: end for

17: // 3) Reshape back to token grid
18:  $G(z) \leftarrow \text{reshape}(h, B \times L \times d)$ 
19: return  $G(z)$ 
```

---

# Mixture of Experts GT

GT-MoE combines explicit geometric inductive biases with gated experts



[Mahadevan, Categories for AGI, 2026]

# Simplicial Message Passing: GT-MoE

Uses a discrete diffusion  
process

Attenuates local  
inconsistencies

Algorithm 10: GT-MoE Mixing Operator  $M$  (Gating + Top- $k$  + Expert Aggregation)

**Input:** Token states  $z \in \mathbb{R}^{B \times L \times d}$ ; experts  $\{E_e\}_{e=1}^E$ ; gating network  $g_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^E$ ; optional top- $k$  value  $k$ ; LayerNorm parameters  $\theta_{\text{LN}}$

**Output:** Mixed token states  $M(z) \in \mathbb{R}^{B \times L \times d}$

1: // 1) Compute gating probabilities per token

2:  $\ell \leftarrow g_\theta(z)$   $\triangleright \ell \in \mathbb{R}^{B \times L \times E}$  (gate logits)

3:  $p \leftarrow \text{softmax}(\ell, \text{ over experts})$   $\triangleright p \in \mathbb{R}^{B \times L \times E}$

4: // 2) Optional sparsification (top- $k$  routing) and renormalization

5: **if**  $k < E$  **then**

6:  $(v, I) \leftarrow \text{TopK}(p, k)$   $\triangleright v, I \in \mathbb{R}^{B \times L \times k}$

7:  $\tilde{p} \leftarrow 0$   $\triangleright$  initialize  $\tilde{p} \in \mathbb{R}^{B \times L \times E}$

8:  $\tilde{p}[\cdot, \cdot, I] \leftarrow v$   $\triangleright$  scatter top- $k$  weights into expert slots

9:  $p \leftarrow \tilde{p} / (\sum_{e=1}^E \tilde{p}_e + \epsilon)$   $\triangleright$  renormalize per token

10: **end if**

11: // 3) Apply experts and aggregate weighted outputs

12:  $u \leftarrow 0 \in \mathbb{R}^{B \times L \times d}$

13: **for**  $e = 1$  **to**  $E$  **do**

14:  $u \leftarrow u + p_e \odot E_e(z)$   $\triangleright p_e$  broadcast to shape  $(B, L, d)$ ;  $\odot$  elementwise

15: **end for**

16: // 4) Residual update and normalization

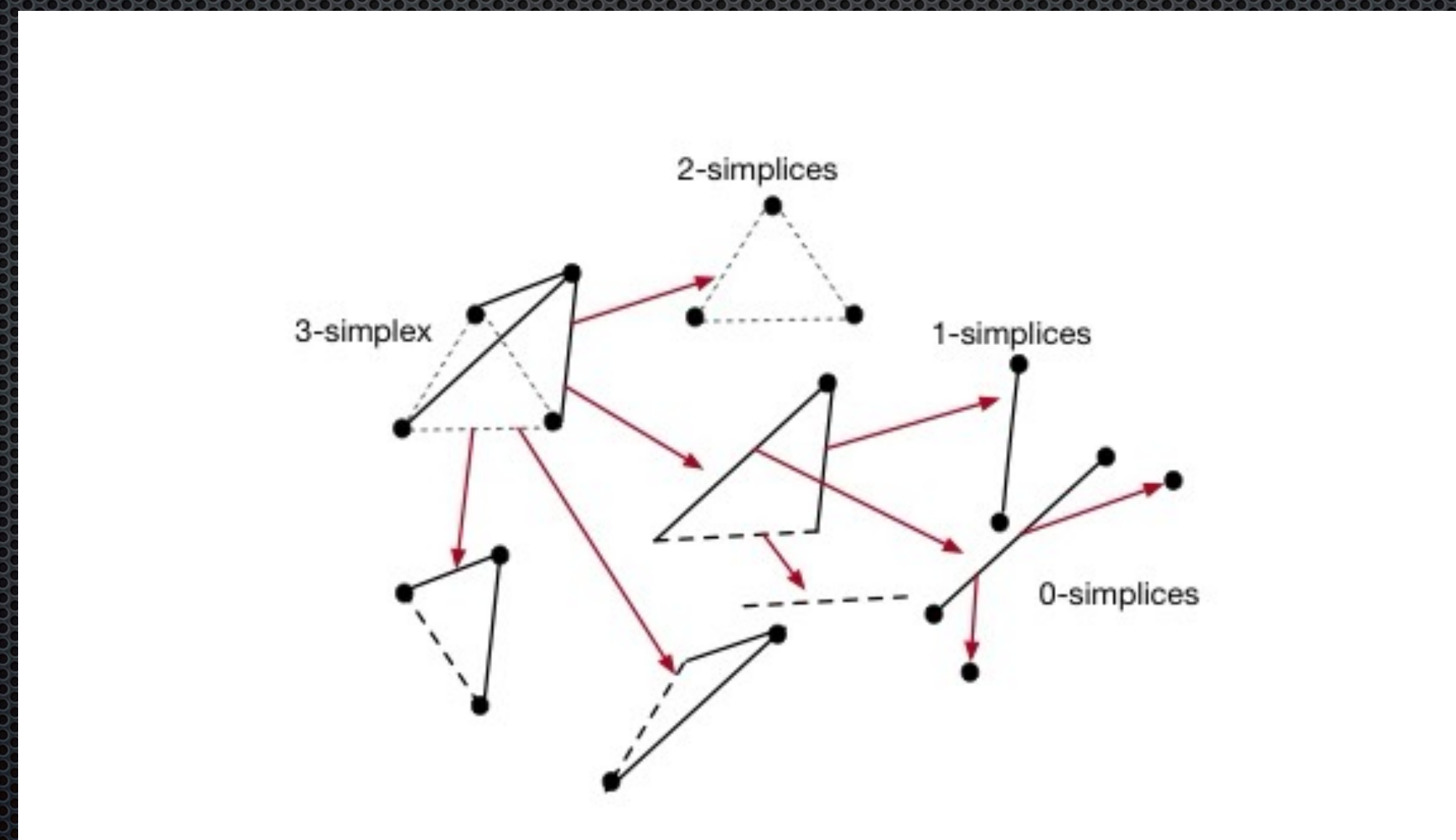
17:  $z' \leftarrow z + u$

18:  $M(z) \leftarrow \text{LayerNorm}_{\theta_{\text{LN}}}(z')$

19: **return**  $M(z)$

# Experimental Results

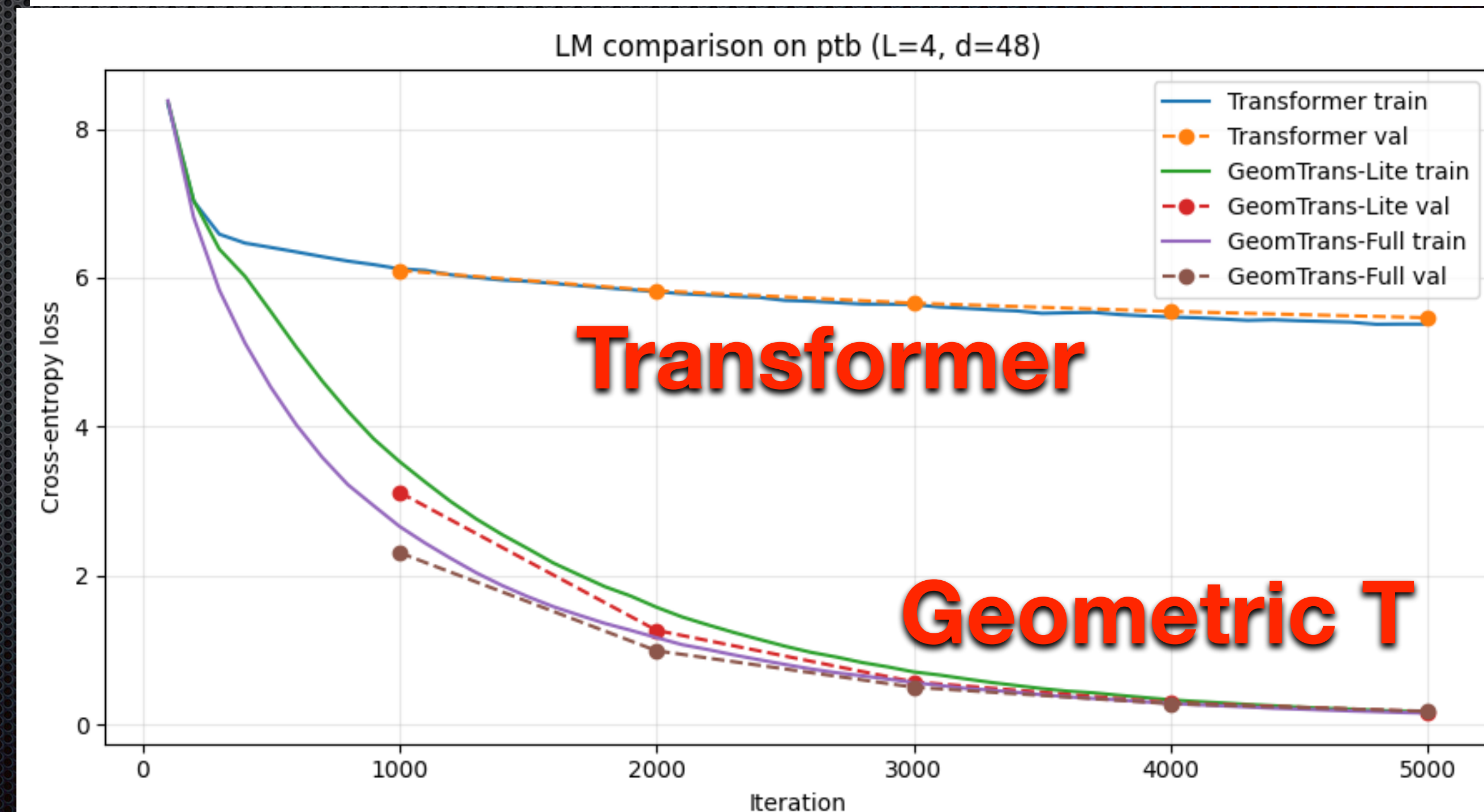
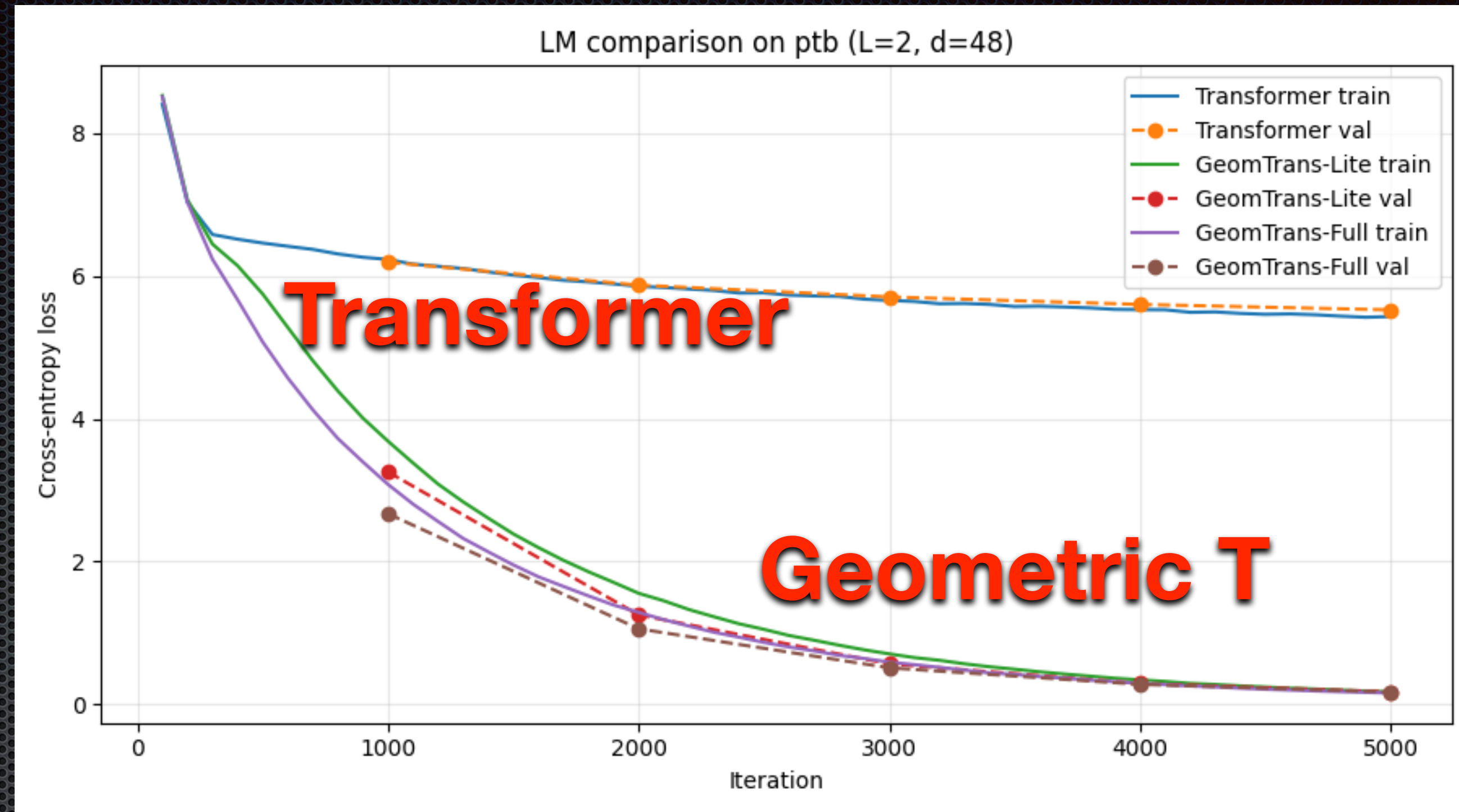
Domain / Task	Base complex (0/1/2-simplices)	Features on 0-simplices	Objective
PTB, WikiText-103 LM	Tokens / adj. pairs / (none)	Word + positional embeddings	Next-token CE loss
Triangle / DB synthetic	Nodes / edges / triangles	Node features (structural)	Triangle / motif classification
FATE causal motifs	Vars / causal edges / motifs	Variable embeddings	Causal motif / FATE score
Large causal manifolds	Vars / causal edges / small motifs	Text embeddings of variables	Unsupervised GT refinement
MNIST diffusion	Pixels / grid edges / $2 \times 2$ patches	Pixel or latent codes	Denosing / diffusion loss



# Language Modeling

## Penn Tree Bank (PTB)

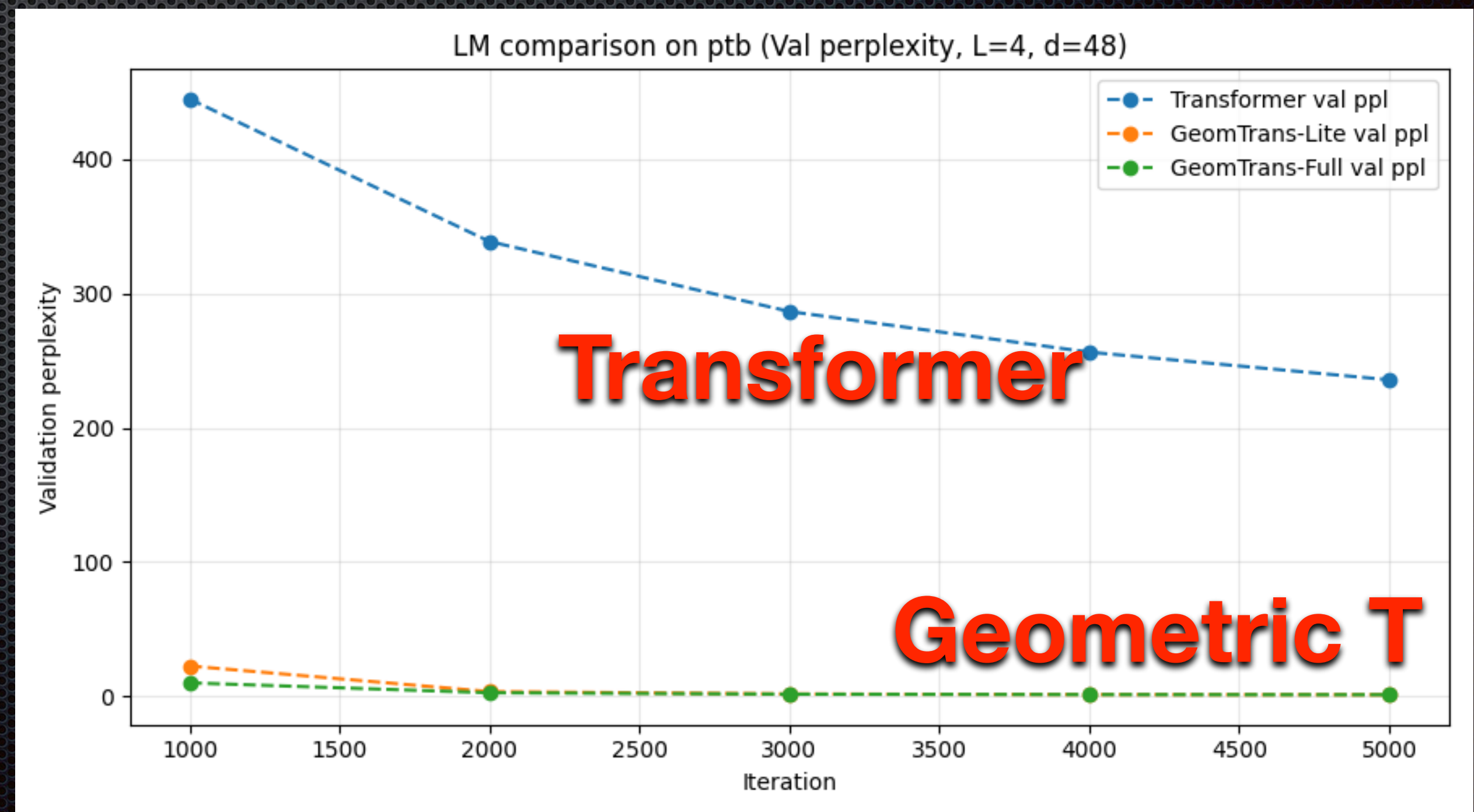
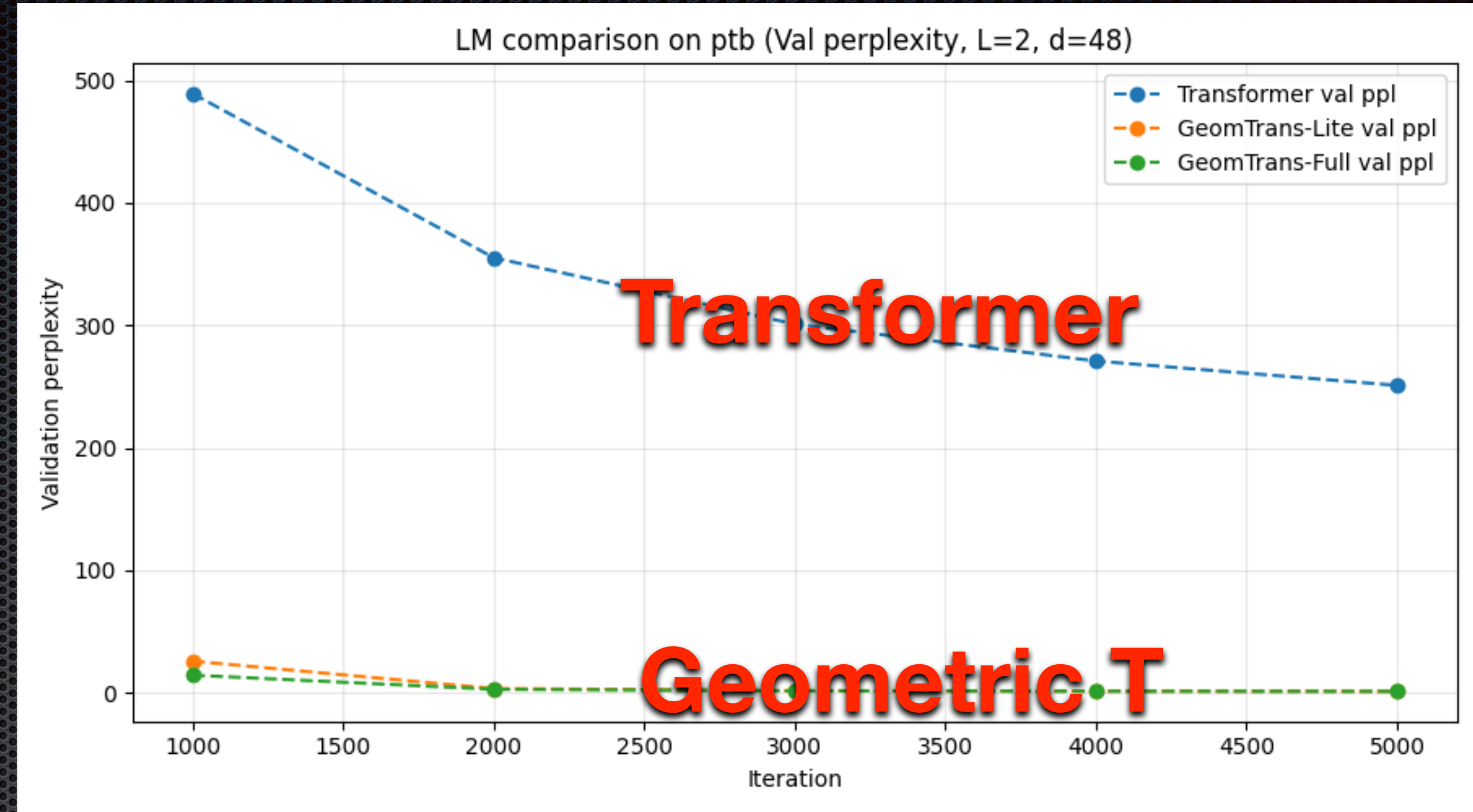
The Penn Treebank (PTB) project selected 2,499 stories from a three year Wall Street Journal (WSJ) collection



# Language Modeling

## Penn Tree Bank (PTB)

The Penn Treebank (PTB) project selected 2,499 stories from a three year Wall Street Journal (WSJ) collection

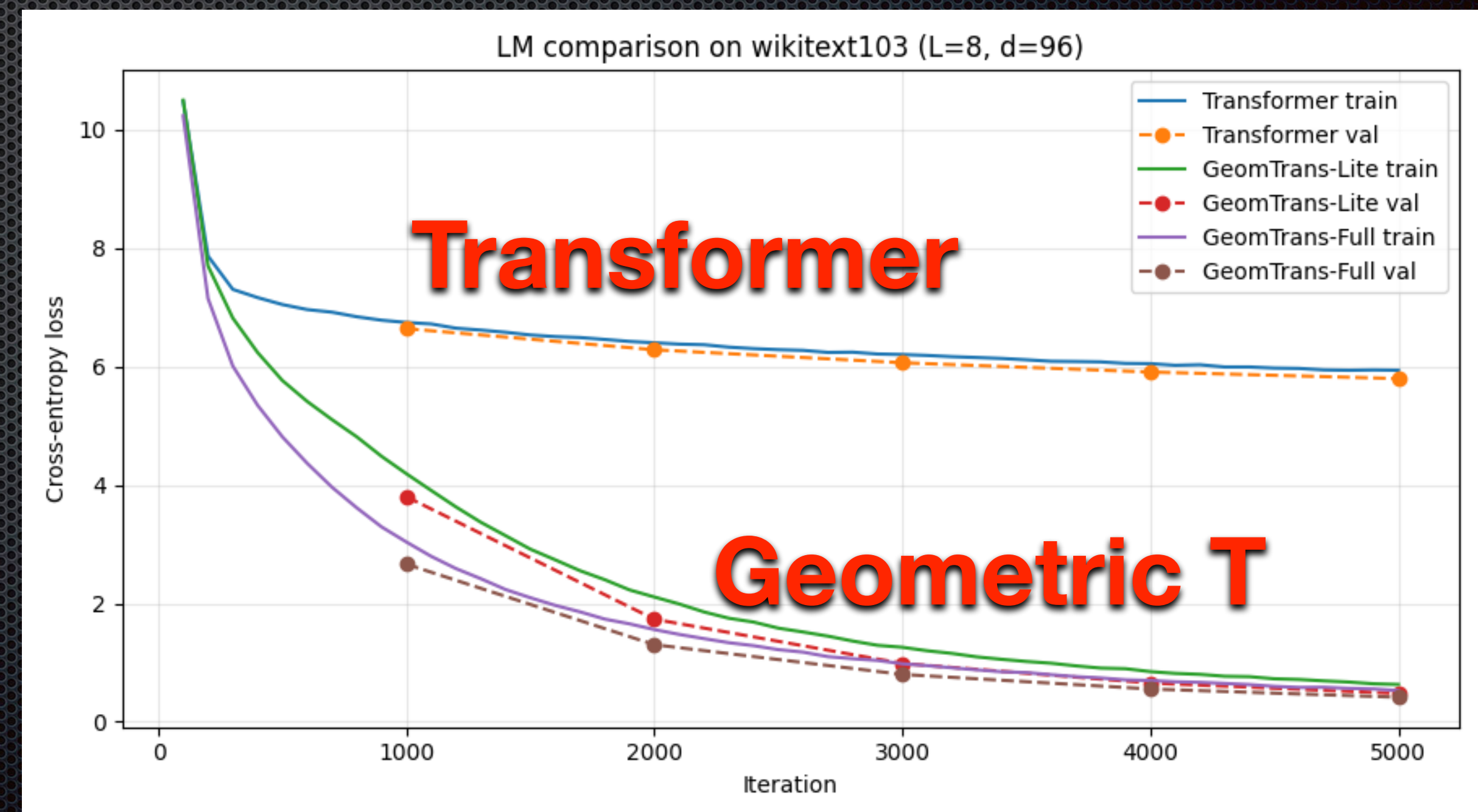
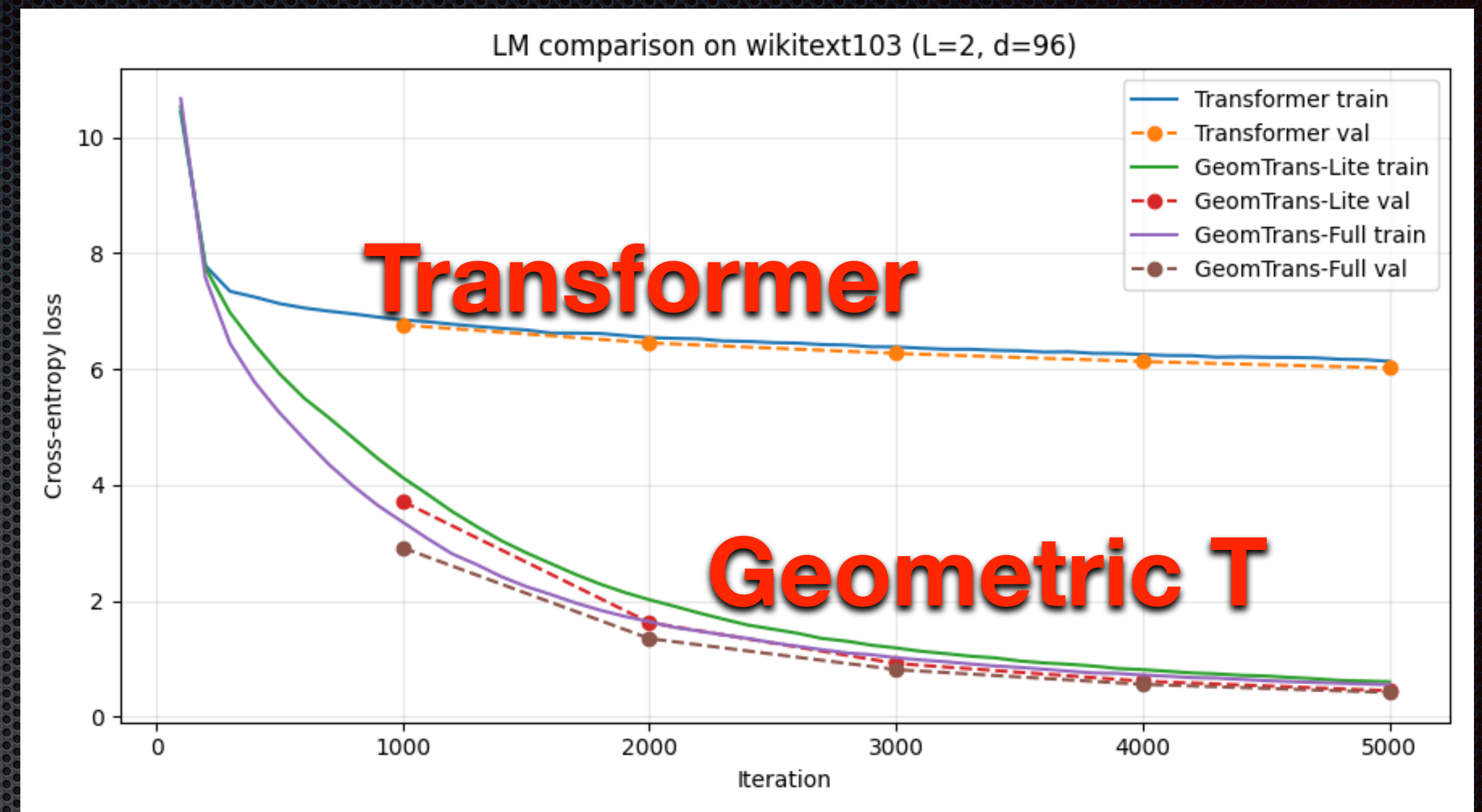


# Language Modeling

Wiki-103: 100-million token dataset of carefully curated Wikipedia articles

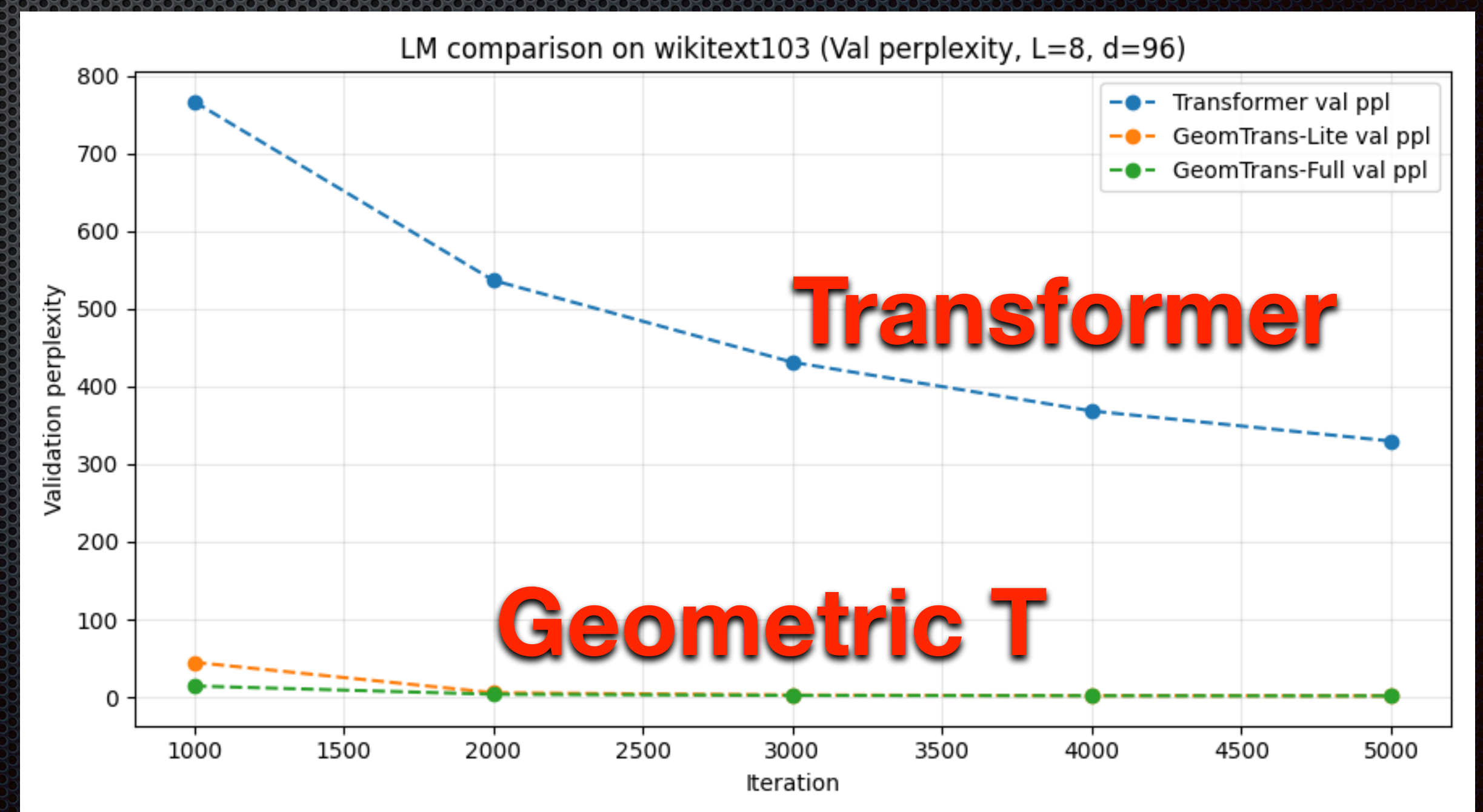
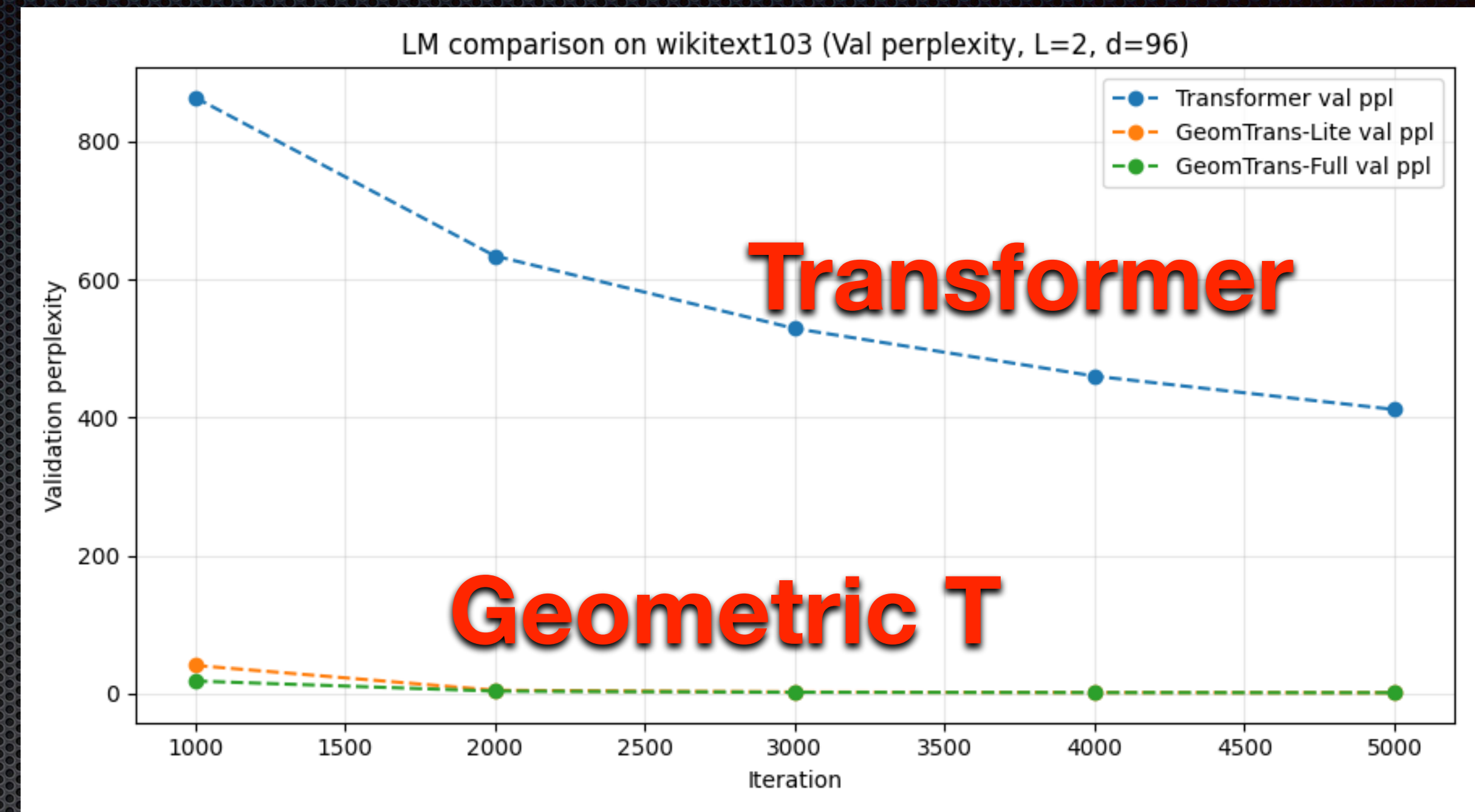
The WikiText-103 dataset include

- 28,475 Wikipedia articles
- Over 100 million words in English
- Complete, untruncated and low-noise texts
- A raw format (TXT), suitable for training autoregressive or bidirectional models



# Language Modeling

Wiki-103: 100-million token dataset of carefully curated Wikipedia articles



# DEMOCRITUS: Causality from Language


[Mahadevan, Large Causal  
Models from Large Language  
Models, Arxiv 2025]

## When did humanity take its first step? Scientists say they now know.

A new analysis of fossils uncovered in Central Africa offers additional evidence that a human ancestor walked upright 7 million years ago.

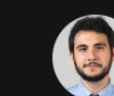
January 2, 2026

🔒 4 min 🌟 Summary ↗ 📌 🗒 289

 Make us preferred on Google



A Sahelanthropus tchadensis skull found in Chad. (Philippe Psaila/Science Source)



By [Dino Grandoni](#)

More than two decades ago, scientists digging in Central Africa unearthed the 7-million-year-old remains of what may be one of the earliest known human ancestors.

## AI Overview

Summary is AI-generated, newsroom-reviewed.

A study published in Science Advances suggests Sahelanthropus tchadensis, a 7-million-year-old ancestor, walked upright, indicating early bipedalism in human evolution. The analysis of limb bones, particularly the femoral tubercle, supports this claim. However, the debate continues as some scientists argue the fossils are too damaged to confirm bipedalism. Further discoveries are needed to resolve the controversy.

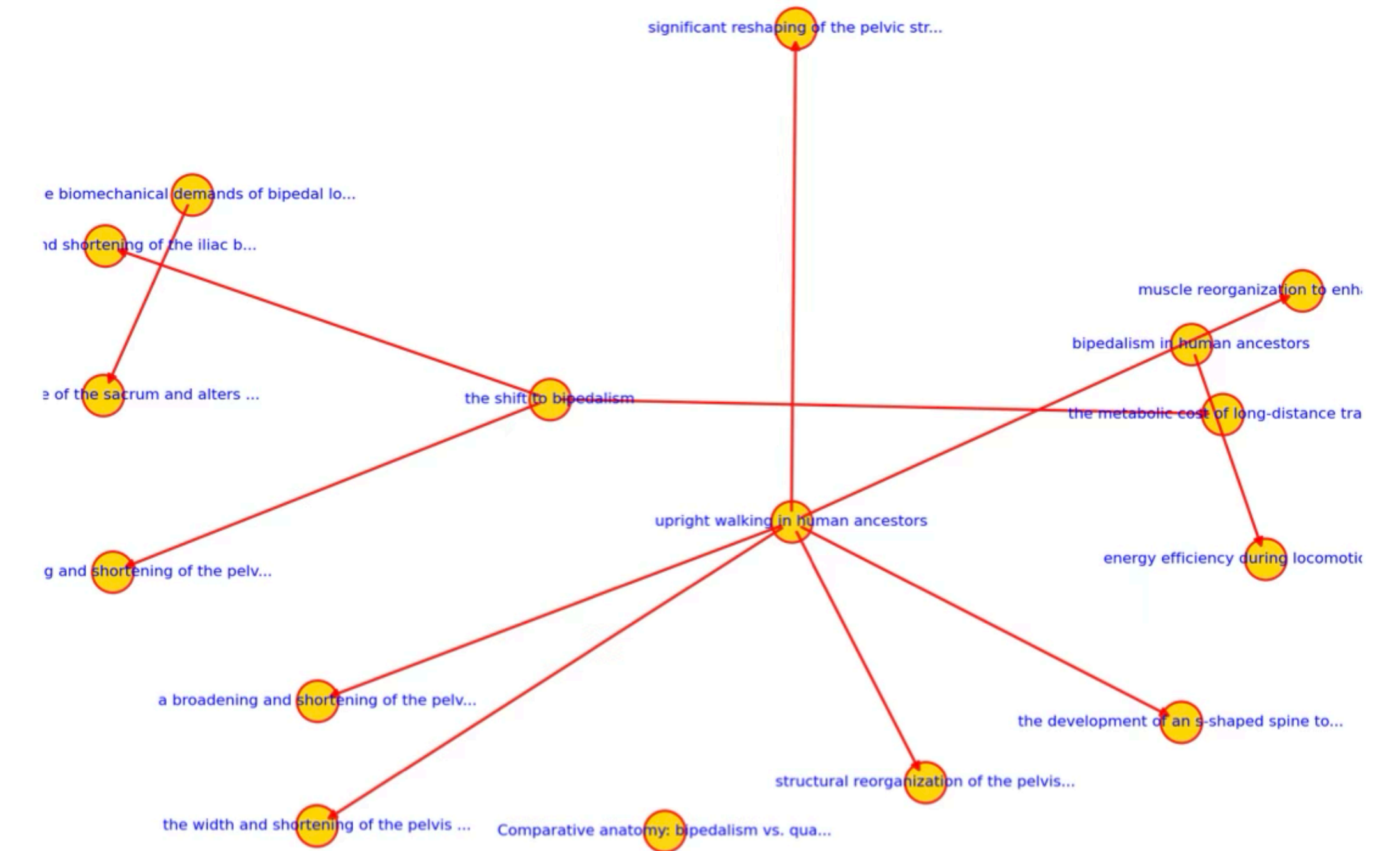
### Read the full article for more on:

- The significance of the femoral tubercle in determining bipedalism.
- Why some scientists remain skeptical about the study's conclusions.
- Future plans for fossil hunting in Chad's Djurab Desert.

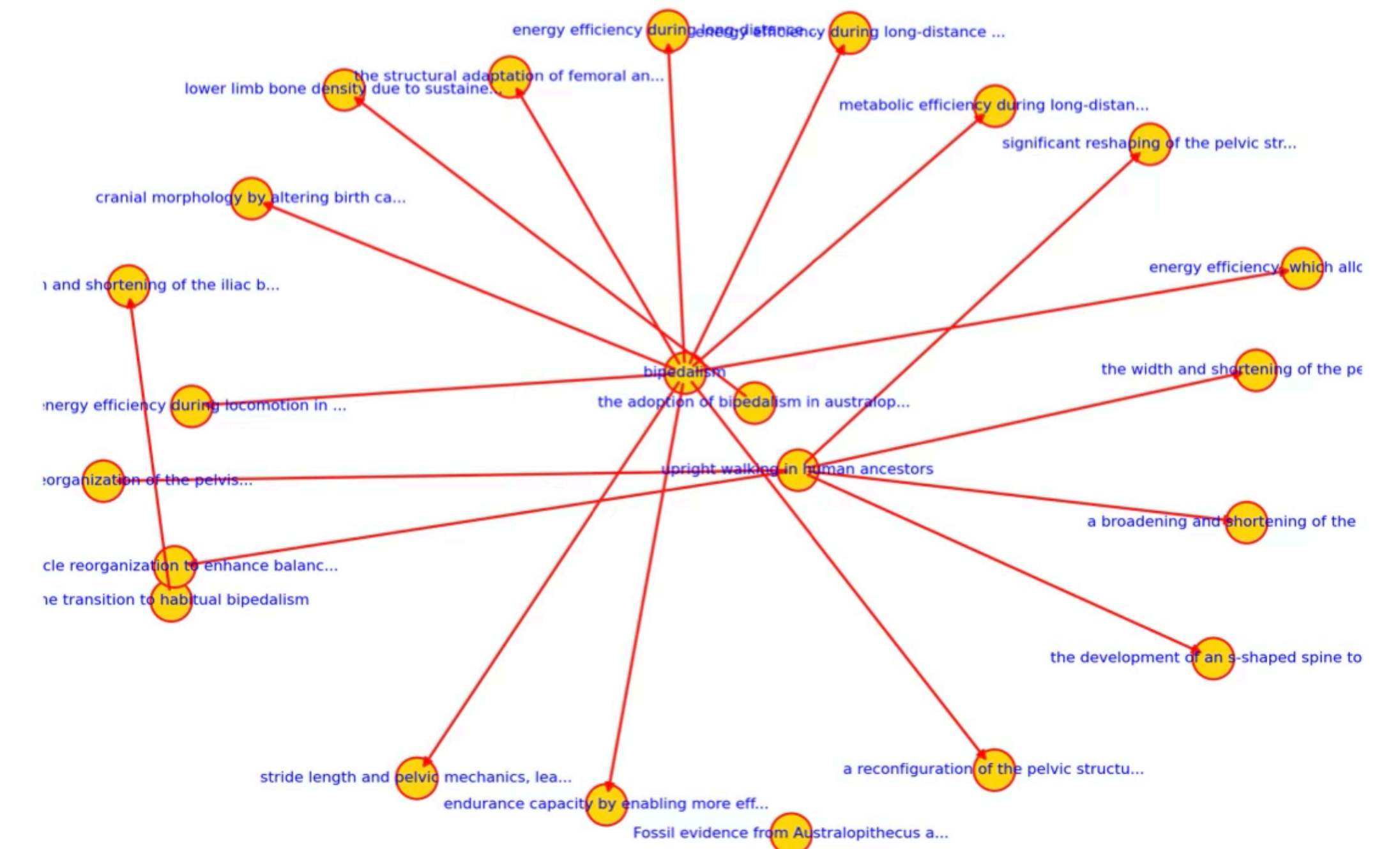
[Did our AI help? Share your thoughts.](#)

# LLM summary

# DEMOCRITUS



manity\_take\_its\_first\_step\_scientists\_say\_they\_no\_70ccf7f26a62 | rank 1 | score=0.155 | Fossil evidence from Australopithecus af...



# Causal DeepDive

DEMOCRITUS explains the causal background to any story

Why, not what!

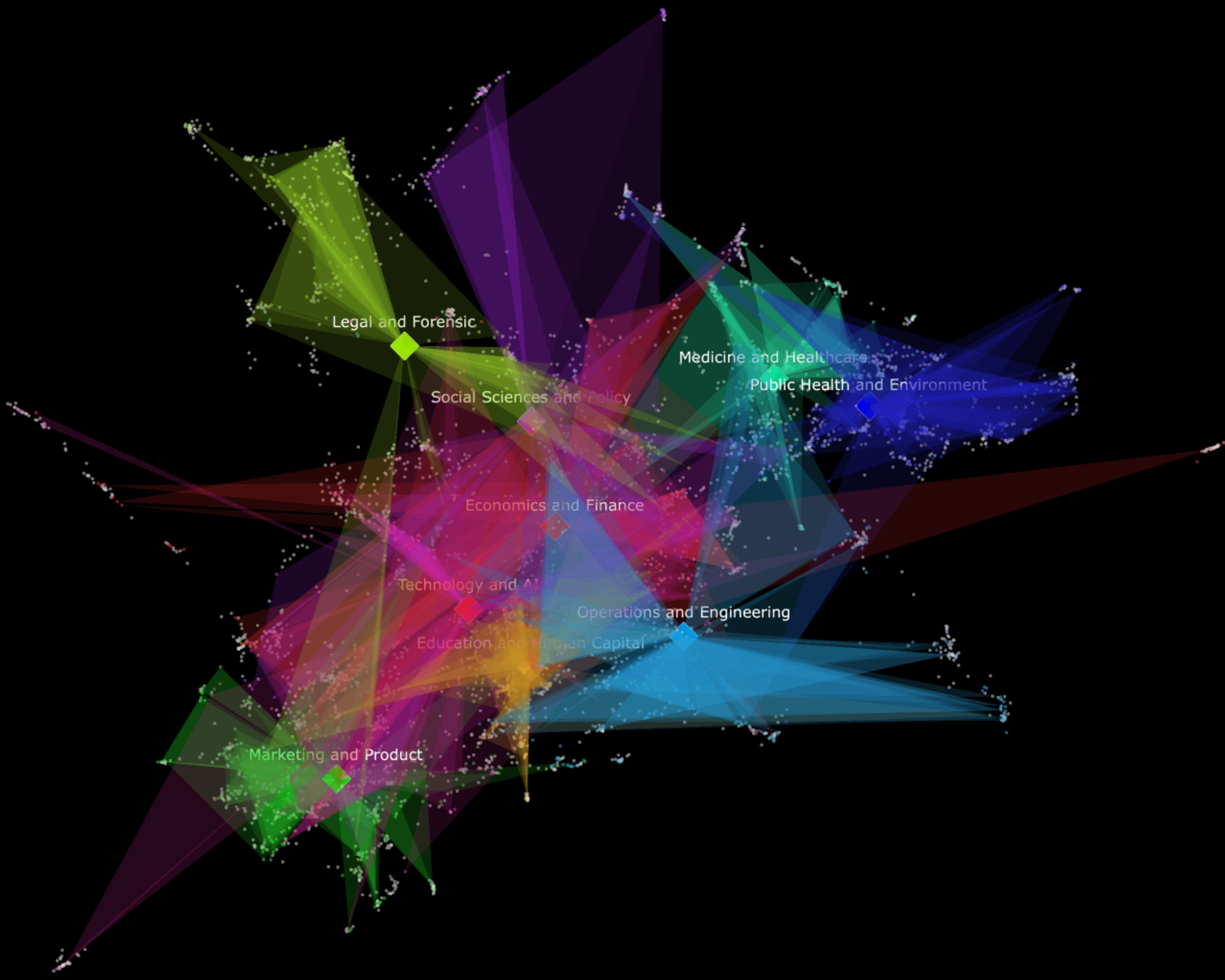
## # Democritus Atlas Summary (v0)

- Nodes: 501
- Edges (unique): 65
- Edge-support rows: 717
- SCC modules (size>1): 0

### ## Top stable bonds (by support\_docs, then score\_sum)

rank	support_docs	support_lcms	score_sum	controversy	rel_type	src	dst
1	1	15	1.462	0.000	INFLUENCES	`bipedalism`	`metabolic efficiency during long-distance travel by optimizing stride mechanics and muscle utilization`
2	1	13	1.328	0.000	INCREASES	`bipedalism`	`energy efficiency during locomotion in early hominins by reducing the metabolic cost of walking over long distances`
3	1	13	1.264	0.000	INCREASES	`bipedalism`	`energy efficiency during long-distance locomotion by reducing the metabolic cost of walking compared to quadrupedal gaits`
4	1	13	1.256	0.000	INFLUENCES	`bipedalism`	`endurance capacity by enabling more effective heat dissipation and sustained pacing over extended distances`
5	1	12	1.209	0.000	INCREASES	`bipedalism`	`energy efficiency which allows for greater metabolic resources to be allocated to brain development`
6	1	11	1.201	0.000	INFLUENCES	`bipedalism`	`structural adaptation of femoral and tibial bones by altering stress distribution patterns over evolutionary time`
7	1	11	1.103	0.000	INFLUENCES	`bipedalism`	`stride length and pelvic mechanics leading to more sustained and efficient terrestrial movement`
8	1	6	0.464	0.000	INFLUENCES	`reduced forest cover`	`selection for upright walking by favoring energy-efficient movement over long distances`
9	1	4	0.403	0.000	INCREASES		`changes in environmental conditions that favored energy-efficient locomotion over long distances`   `likelihood of bipedalism in early hominins`
10	1	6	0.156	0.000	INFLUENCES	`abnormal femoral tubercle position`	`knee joint alignment during locomotion by altering the line of pull on the patellar tendon`
11	1	2	0.136	0.000	INFLUENCES	`fossil structure of ardpithecus`	`emergence of efficient terrestrial walking in early hominins by demonstrating adaptations in the pelvis and foot that support bipedal locomotion`
12	1	6	0.129	0.000	CAUSES	`abnormal femoral tubercle position`	`increased postoperative patellar instability after tibial tubercle osteotomy`
13	1	2	0.128	0.000	INCREASES		`muscle reorganization for balance and propulsion in bipedalism`   `efficiency of upright walking in human ancestors`
14	1	6	0.081	0.000	CAUSES	`larger brain size in primates`	`enhanced problem-solving abilities by enabling greater neural complexity and cognitive flexibility`
15	1	4	0.076	0.000	INCREASES	`lateral displacement of the femoral tubercle`	`patellar maltracking leading to altered patellofemoral alignment`

# 100,000 causal claims



Domain	# Triangles
Economics and Finance	1140
Education and Human Capital	950
Legal and Forensic	553
Marketing and Product	1059
Medicine and Healthcare	932
Operations and Engineering	1281
Public Health and Environment	1336
Social Sciences and Policy	1015
Technology and AI	782

*Table 1.* Simplicial complex induced by local causal neighborhoods in the 100K-claim dataset.

# Summary

- ✦ The ordinal category  $\Delta$  has as objects  $[n]$ , and arrows weakly order-preserving functions on  $[n]$
- ✦ Simplicial objects are contravariant functors  $X[n] : \Delta \rightarrow \mathcal{C}$
- ✦ GAIA: Generative AI Architecture based on simplicial objects
- ✦ Geometric Transformer is a GAIA architecture
- ✦ Diagrammatic Backpropagation is used to train a GT model