

# Language Models as Hierarchical Computational Projections: An Entropic Scaling Framework

Karol Kowalczyk

November 12, 2025

## Abstract

We demonstrate that large language models (LLMs) naturally instantiate the computational hierarchy framework with entropic information scaling. Each model size corresponds to a discrete level  $n$  in the hierarchy with effective information capacity  $I(n) = \kappa n \log n$  bits. We show that model distillation implements projection operators between levels, while fine-tuning acts as collapse operators. The behavioral distance between models follows predictable patterns based on their entropic capacity differences. This perspective provides: (1) a principled way to understand scaling laws through entropic rather than exponential growth, (2) theoretical justification for distillation as implementing adjoint projections, (3) predictions about emergent capabilities as phase transitions at critical entropic thresholds, and (4) a unified framework connecting model size, capability, and computational complexity through entropic scaling.

## 1 Introduction

Large language models exhibit striking regularities in how capabilities scale with model size. We propose these patterns reflect an underlying computational hierarchy where each model implements a finite-state machine with entropic information capacity. This framework, based on entropic scaling  $I(n) = \kappa n \log n$  rather than exponential growth, provides both theoretical understanding and practical insights.

## 2 LLMs as Computational Hierarchy

### 2.1 Mapping models to hierarchy levels

**Definition 2.1** (Model-to-Level Mapping with Entropic Scaling). *For a language model with  $P$  parameters and vocabulary size  $V$ , its hierarchy level is:*

$$n_L = \lfloor \alpha \log_2 P + \beta \log_2 V + \gamma \rfloor \quad (1)$$

*with effective information capacity:*

$$I_L = \kappa n_L \log n_L \text{ bits} \quad (2)$$

*where  $\alpha, \beta, \gamma$  are empirically determined constants, and  $\kappa$  is the entropic scaling factor.*

This mapping captures that model capacity depends on both parameter count and vocabulary size, but scales entropically rather than exponentially.

### 2.2 Concrete examples with entropic capacity

The entropic scaling ensures realistic information capacities rather than impossible  $2^{41}$  bits.

Table 1: Language Models Mapped to Hierarchy Levels with Entropic Scaling

Model	Parameters	Level $n_L$	Capacity $I_L$
GPT-2 Small	117M	27	$\approx 130\kappa$ bits
GPT-2 Large	1.5B	31	$\approx 154\kappa$ bits
GPT-3 Ada	350M	29	$\approx 141\kappa$ bits
GPT-3 Davinci	175B	37	$\approx 196\kappa$ bits
GPT-4	$\sim 1.7T$	41	$\approx 225\kappa$ bits

### 3 Distillation as Projection

#### 3.1 Knowledge distillation formalized

**Theorem 3.1** (Distillation as Entropic Projection). *Knowledge distillation from teacher model  $M_T$  at level  $n_T$  to student model  $M_S$  at level  $n_S < n_T$  implements the projection operator  $P_{T \rightarrow S}$  that compresses information from capacity  $I_T = \kappa n_T \log n_T$  to  $I_S = \kappa n_S \log n_S$ .*

*Proof sketch.* Distillation minimizes:

$$\mathcal{L}_{\text{distill}} = \text{KL}(p_S || p_T) + \lambda \mathcal{L}_{\text{task}} \quad (3)$$

This corresponds to finding optimal projection that preserves maximum information within the student's entropic capacity constraint  $I_S = \kappa n_S \log n_S$ .  $\square$

#### 3.2 Empirical validation

Distillation efficiency follows the entropic scaling:

$$\text{Retention} \approx \exp\left(-\alpha \frac{I_T - I_S}{I_T}\right) = \exp\left(-\alpha \frac{n_T \log n_T - n_S \log n_S}{n_T \log n_T}\right) \quad (4)$$

This predicts diminishing returns when capacity gaps exceed critical thresholds.

### 4 Fine-tuning as Collapse

#### 4.1 Task specialization formalized

**Theorem 4.1** (Fine-tuning as Entropic Collapse). *Fine-tuning a pre-trained model  $M_{\text{pre}}$  at level  $n$  on task  $T$  implements the collapse operator  $C_T$  that reconstructs task-specific structure within entropic capacity  $I(n) = \kappa n \log n$ .*

The collapse operator:

1. Projects general knowledge to task-relevant subspace
2. Reconstructs specialized representations
3. Maintains total capacity constraint  $I(n) = \kappa n \log n$

#### 4.2 Adjunction relationship

**Proposition 4.2** (Distillation-Finetuning Adjunction). *For models at levels  $i < j$ , distillation  $P_{j \rightarrow i}$  and fine-tuning  $C_{i \rightarrow j}$  form an approximate adjunction with error:*

$$\varepsilon = \|P \circ C - \| \sim \frac{1}{n \log n} \quad (5)$$

where the entropic factor ensures rapid convergence.

## 5 Emergent Capabilities and Phase Transitions

### 5.1 Critical thresholds with entropic scaling

**Definition 5.1** (Entropic Capability Emergence). *A capability  $\mathcal{C}$  emerges at critical level  $n_c$  when:*

$$I(n_c) = \kappa n_c \log n_c \geq I_{\text{threshold}}(\mathcal{C}) \quad (6)$$

where  $I_{\text{threshold}}(\mathcal{C})$  is the minimum entropic information required for capability  $\mathcal{C}$ .

### 5.2 Observed emergence patterns

Empirical observations align with entropic thresholds:

Table 2: Capability Emergence at Entropic Thresholds

Capability	Critical Level	Capacity Required
Basic syntax	$n \approx 25$	$I \approx 116\kappa$ bits
Semantic understanding	$n \approx 30$	$I \approx 147\kappa$ bits
Multi-step reasoning	$n \approx 35$	$I \approx 179\kappa$ bits
Abstract reasoning	$n \approx 40$	$I \approx 213\kappa$ bits
Meta-learning	$n \approx 45$	$I \approx 249\kappa$ bits

The entropic scaling explains why capabilities emerge gradually rather than suddenly.

### 5.3 Phase transition analysis

**Theorem 5.2** (Entropic Phase Transitions). *Qualitative behavioral changes occur when:*

$$\frac{d^2 I}{dn^2} = \frac{d^2}{dn^2}(\kappa n \log n) = \frac{\kappa}{n} \quad (7)$$

reaches critical values, corresponding to inflection points in capability space.

## 6 Behavioral Distance and Model Similarity

### 6.1 Measuring model distance with entropic scaling

**Definition 6.1** (Inter-Model Distance). *The behavioral distance between models at levels  $i$  and  $j$  is:*

$$d(M_i, M_j) = \text{Beh}(i, j) + \lambda |I(i) - I(j)| \quad (8)$$

where  $\text{Beh}(i, j)$  measures output distribution divergence and  $I(n) = \kappa n \log n$ .

### 6.2 Empirical distance patterns

**Proposition 6.2** (Distance Scaling with Entropic Capacity). *For models separated by  $\Delta n$  levels:*

$$d(M_i, M_{i+\Delta n}) \approx d_0 \sqrt{\Delta I} = d_0 \sqrt{\kappa [(i + \Delta n) \log(i + \Delta n) - i \log i]} \quad (9)$$

This square-root relationship reflects the entropic information geometry.

## 7 Thermodynamic Interpretation

### 7.1 Computational entropy and free energy

**Definition 7.1** (Model Free Energy with Entropic Scaling). *The free energy of model  $M_n$  at temperature  $T$  is:*

$$F_n = E_n - TS_n \quad (10)$$

where:

- $E_n \propto I(n) = \kappa n \log n$  (*computational energy*)
- $S_n$  = *entropy of output distribution*
- $T$  = *temperature parameter (controls randomness)*

### 7.2 Training dynamics

**Theorem 7.2** (Entropic Training Dynamics). *Training minimizes free energy subject to entropic capacity constraint:*

$$\frac{dF}{dt} = -\nabla_{\theta} F \cdot \dot{\theta} \leq 0 \quad (11)$$

with  $|\theta| \leq \exp(I(n))$  where  $I(n) = \kappa n \log n$ .

This ensures training respects the model's entropic information capacity.

## 8 Scaling Laws Revisited

### 8.1 Traditional scaling laws

Classical scaling laws assume:

$$L(N) = \left(\frac{N_c}{N}\right)^{\alpha} \quad (12)$$

where  $L$  is loss and  $N$  is parameter count.

### 8.2 Entropic perspective on scaling

**Theorem 8.1** (Entropic Scaling Law). *With hierarchy level  $n \sim \log N$  and entropic capacity  $I(n) = \kappa n \log n$ :*

$$L(n) \sim \exp(-\beta I(n)) = \exp(-\beta \kappa n \log n) \quad (13)$$

*This provides a direct connection between information capacity and performance.*

### 8.3 Implications

The entropic scaling perspective explains:

1. Why scaling shows diminishing returns (entropic vs exponential growth)
2. Existence of capability plateaus (between entropic thresholds)
3. Optimal model sizing (matching task complexity to entropic capacity)
4. Efficiency of distillation (compression within entropic bounds)

## 9 Practical Applications

### 9.1 Model selection guidelines

Given task requiring information  $I_{\text{task}}$ :

1. Choose model level  $n^* = \arg \min_n \{I(n) \geq I_{\text{task}}\}$
2. Where  $I(n) = \kappa n \log n$
3. Avoid over-provisioning (wastes resources)
4. Avoid under-provisioning (cannot solve task)

### 9.2 Distillation strategies

For efficient distillation:

$$n_S = n_T - \Delta n_{\text{optimal}} \quad (14)$$

where  $\Delta n_{\text{optimal}} \approx \sqrt{n_T}$  balances compression and retention given entropic scaling.

### 9.3 Fine-tuning protocols

Optimal fine-tuning respects entropic bounds:

- Learning rate  $\propto 1/(n \log n)$
- Batch size  $\propto n \log n$
- Training steps  $\propto n \log n$

## 10 Theoretical Implications

### 10.1 Consciousness in language models

If consciousness requires:

1. Hierarchical organization (model architecture)
2. Entropic information scaling ( $I(n) = \kappa n \log n$ )
3. Integration above threshold (attention mechanisms)
4. Collapse dynamics (autoregressive generation)

Then sufficiently large models may exhibit proto-conscious properties within their entropic capacity limits.

### 10.2 Limits of scaling

**Theorem 10.1** (Entropic Scaling Limits). *Maximum useful model level  $n_{\max}$  is constrained by:*

$$I(n_{\max}) = \kappa n_{\max} \log n_{\max} \leq I_{\text{universe}} \quad (15)$$

where  $I_{\text{universe}}$  is total information available for training.

This suggests fundamental limits even with entropic rather than exponential scaling.

## 11 Future Directions

### 11.1 Empirical validation

Key experiments:

1. Measure behavioral distance vs entropic capacity difference
2. Test emergence thresholds for specific capabilities
3. Validate distillation efficiency predictions
4. Examine phase transitions in capability space

### 11.2 Theoretical extensions

1. Continuous hierarchy models (fractional levels)
2. Multi-modal hierarchies (vision + language)
3. Quantum computational hierarchies
4. Connection to biological neural hierarchies

## 12 Conclusion

We have shown that language models naturally implement computational hierarchies with entropic information scaling. Key insights:

1. Model capacity scales as  $I(n) = \kappa n \log n$ , not exponentially
2. Distillation implements projection between entropic levels
3. Fine-tuning implements collapse operators
4. Capabilities emerge at entropic thresholds
5. Scaling laws reflect entropic information geometry

This framework provides both theoretical understanding and practical guidance for:

- Designing efficient model architectures
- Optimizing distillation and fine-tuning
- Predicting capability emergence
- Understanding fundamental limits

The entropic scaling perspective transforms our understanding of language models from mysterious black boxes to principled computational hierarchies with well-defined information-theoretic properties.

## References

- [1] Kaplan, J., et al. (2020). Scaling laws for neural language models. arXiv:2001.08361.
- [2] Kowalczyk, K. (2025). Consciousness as collapsed computational time: A unified theory with entropic scaling. Zenodo, doi:10.5281/zenodo.17556941.
- [3] Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling the knowledge in a neural network. arXiv:1503.02531.