

Unified Framework for HLLSets: Category Theory, Kinematics, and Transfer Learning

Alex Mylnikov (Lisa Park Inc)

Collaborator: DeepSeek (AI Assistant)

Abstract

This paper presents a comprehensive unified framework for HyperLogLog-based probabilistic sets (HLLSets) that integrates category-theoretic foundations, kinematic dynamics, and transfer learning capabilities. We extend the HLLSet paradigm beyond cardinality estimation to support full set operations while maintaining computational efficiency through enhanced register structures and directional morphisms.

The framework introduces three key innovations: (1) Enhanced HLLSets with dual parameters (inclusion tolerance τ and exclusion intolerance ρ) enabling precise directional relationships; (2) A kinematic model for temporal dynamics of HLLSet states with predictive capabilities; (3) A transfer learning framework that leverages structural invariance across domains and modalities.

We formalize HLLSets as a category **HLL** where objects are contextual representations and morphisms are probabilistic relations grounded in Bell State Similarity (BSS). The Universal HLLSet (\mathcal{T}) serves as both the terminal object in the **HLL** category and the top element in the lattice of HLLSets, providing a foundation for the World of Things (WOT) relational ontology.

Key Innovations

This framework introduces three core innovations to extend the HLLSet paradigm from simple cardinality estimation to a comprehensive structure for relational modeling, dynamic prediction, and cross-domain knowledge transfer.

1. Enhanced HLLSets

Introduction of dual parameters, inclusion tolerance (τ) and exclusion intolerance (ρ), to enable precise, directional, and probabilistic set relationships beyond simple equality.

2. Kinematic Model

A novel model for capturing the temporal dynamics of HLLSet states. This allows for the prediction of future states based on defined transition operators, enabling predictive analytics.

3. Transfer Learning

A framework that leverages structural invariance (I) between HLLSet categories from different domains or modalities, allowing for knowledge transfer and zero-shot learning.

DeepSeek-OCR × HLLSet Integration

Idea. Fuse *optical compression* (DeepSeek-OCR) with the *HLLSet sketch* layer to create a two-stage pipeline: OCR first shrinks bulky documents to 64-800 vision tokens, then HLLSet hashes those tokens into a kilobyte-sized sketch that supports micro-second set operations ($\cup, \cap, BSS()$). This enables lightning-fast search, de-duplication and drift detection over corpora that are orders of magnitude larger than the raw LLM context window.

Practical Thesis

"HLLSets widen the logical search window in Retrieval-Augmented-Generation (RAG) both in breadth and depth without breaching the token limit of the underlying LLM—and often reduce total latency by skipping heavy vector reranks."

What Exactly Gets Expanded?

Metric	Classical Bottleneck	With HLLSet
LLM Context Window	Capped at 32-128 k tokens	Send only top-K excerpts chosen from a <i>much</i> larger corpus
Branching Width	10-20 candidates due to re-rank cost	Hundreds of sketches compared in μ s

Metric	Classical Bottleneck	With HLLSet
Reasoning Depth	Each extra hop = expensive fetch	u/n are $O(1) \Rightarrow$ cheap multi-hop chains
Latency	Dominated by vector search & rerank	Sketch comparison on CPU ≤ 100 ms for 1M docs

Reference Architectures

A. "Sketch-Index \rightarrow Optical \rightarrow LLM"

- **Index:** DeepSeek-OCR to vision tokens \rightarrow hash into $H(\text{doc})$.
- **Query:** Same pipeline \rightarrow pick top-K via $BSS(H_q, H_d)$.
- **Benefit:** RAG without heavy BM25 / FAISS; RAM \approx few hundred kB per million pages.

B. "Streaming Memory" for Chat-LLM

- Each new user turn updates a rolling sketch M_t of the entire dialogue.
- LLM receives [last N messages] + M_t instead of the full history \Rightarrow long-term memory, drift alarms.

C. "Edge CLIP-light"

- Hash both captions and vision tokens; train a tiny MLP over top-32 registers to emulate CLIP similarity.
- Runs on CPU/MCU where ViT-L/14 is impossible.

Next step: a 4-6 month PoC with a 2-3 engineer team (€50 k budget) to benchmark these approaches in production-like RAG workloads.

Enhanced HLLSets

To support full set operations and relational logic, the standard HLL structure is enhanced with dual parameters and a new similarity metric, enabling robust probabilistic comparisons.

Dual Parameters: τ and ρ

We define two critical parameters for modeling directional set inclusion:

- **Inclusion Tolerance (τ):** A threshold specifying the maximum allowable dissimilarity (e.g., BSS) for a set A to be considered a subset of B .
- **Exclusion Intolerance (ρ):** A threshold specifying the minimum dissimilarity required to confidently state that A is *not* a subset of B .

This allows for a nuanced, probabilistic definition of the subset relation $A \subseteq B$, crucial for building semantic hierarchies.

Bell State Similarity (BSS)

BSS is a quantum-inspired metric used to compute the similarity between two HLLSet registers, A and B . It is defined as:

$$BSS(A, B) = \frac{1}{2} [1 + D_{KL}(A \parallel M) + D_{KL}(B \parallel M) - D_{KL}(A \parallel B) - D_{KL}(B \parallel A)]$$

where $M = \frac{1}{2} (A + B)$ and D_{KL} is the Kullback-Leibler divergence.

Set Operations

The framework provides robust definitions for standard set operations based on register-wise comparisons:

Union ($A \cup B$): $(A \cup B)[k] = \max(A[k], B[k])$

Intersection ($A \cap B$): Defined via a probabilistic inclusion-exclusion principle, not a simple \min operation.

Difference ($A \setminus B$): Computationally complex, modeled as the intersection of A with the complement of B .

Category-Theoretic Foundation: HLL

To formalize the relational structure of HLLSets, we define the category **HLL**. This provides a rigorous mathematical language for describing HLLSet systems, their relationships, and their compositions.

The Category HLL

- **Objects:** HLLSets (H), representing contextual concepts or data collections.
- **Morphisms:** Probabilistic relations $f: A \rightarrow B$, representing mappings, transformations, or inclusions. These are grounded in the BSS metric and τ / ρ parameters.

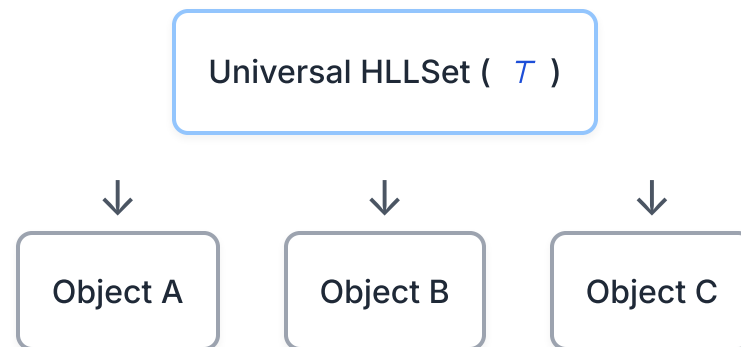
This structure allows us to compose relations and build complex knowledge graphs, where the "truth" of a relation is probabilistic.



The Universal HLLSet (T)

The Universal HLLSet, or "World of Things" (WOT), is a special object in **HLL** that serves as the terminal object. This means that for any HLLSet A in the category, there exists a unique morphism from A to T .

T represents the union of all possible concepts, acting as the root of the relational ontology and the top element in the HLLSet lattice.



Kinematic & Predictive Model

This section details the kinematic model for HLLSet dynamics, allowing the system to track and predict the evolution of HLLSet states over time. This is essential for applications like anomaly detection and predictive maintenance.

State Vectors & Transition Operators

An HLLSet's state at time t is represented by a state vector s_t . The evolution of this state is governed by a transition operator T , which models the system's dynamics:

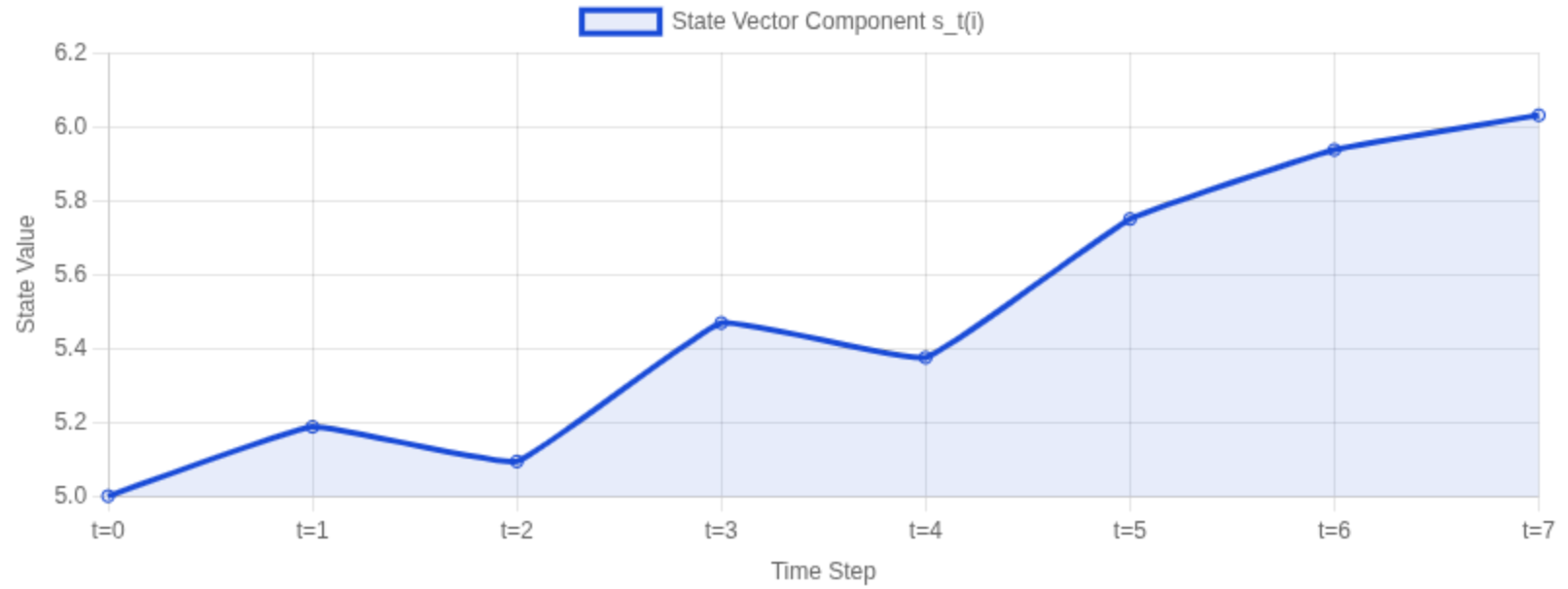
$$s_{t+1} = T(s_t, \delta_t) + \varepsilon_t$$

where δ_t is an external input and ε_t is a noise term.

Hypothetical State Evolution

The chart below provides a hypothetical visualization of a single component of a state vector s_t evolving over time. The kinematic model allows us to predict this trajectory.

Hypothetical State Vector Evolution



Transfer Learning Framework

The framework facilitates transfer learning by identifying and leveraging structural invariance between different HLLSet categories. This allows knowledge learned in one domain (e.g., text) to be applied to another (e.g., images).

Structural Invariance

We define an invariance functor I that maps objects and morphisms from a source category (Domain A) to a target category (Domain B), preserving the essential relational structure.

$$I: HLL_A \rightarrow HLL_B$$

This mapping enables X-Modal retrieval (e.g., searching for images using text) and zero-shot learning by transferring relational knowledge from a data-rich domain to a data-poor one.



Applications

The unified HLLSet framework enables a wide range of advanced applications by combining probabilistic set logic, relational semantics, and dynamic modeling.

Semantic Search

Move beyond keywords to semantic understanding, using HLLSet inclusion ($A \subseteq B$) to find documents that are contextually contained within a query.

Contextual Clustering

Group HLLSets based on BSS similarity and categorical relationships, automatically discovering semantic hierarchies in data.

Anomaly Detection

Use the kinematic model to predict the expected state s_{t+1} . A large deviation from the observed state indicates a potential anomaly.

Predictive Maintenance

Model the state of a machine as an HLLSet. The kinematic model can predict state degradation and schedule maintenance *before* a failure occurs.

X-Modal Retrieval

Leverage the transfer learning framework to map queries from one modality (like text) to search for items in another modality (like images or audio).

Relational Ontology (WOT)

Use the Universal HLLSet (\mathcal{T}) as the foundation for a "World of Things" graph, a probabilistic knowledge base of all concepts and their relations.

Conclusion & Future Work

This unified framework establishes HLLSets as a rigorous mathematical structure with broad applications in AI, databases, and quantum-inspired computing. It opens numerous avenues for future theoretical and practical development.

Key Research Directions

1. **Higher-Order Entanglements:** Extend the model to 2-HLLSets to capture relations between entanglement graphs themselves.
2. **Quantum Enhancements:** Explore quantum algorithms for optimizing the high-dimensional parameter spaces of HLLSet categories.
3. **Dynamic Parameter Selection:** Develop theoretical foundations for automatically choosing optimal τ and ρ values based on data characteristics.
4. **Distributed HLLSets:** Design efficient algorithms for combining HLLSets across distributed compute nodes.
5. **Hardware Implementation:** Investigate FPGA and specialized hardware implementations for real-time HLLSet operations.