# LITS: An Optimized Learned Index for Strings

Yifan Yang, Shimin Chen

SKLP, ACS, Institute of Computing Technology, CAS

University of Chinese Academy of Sciences
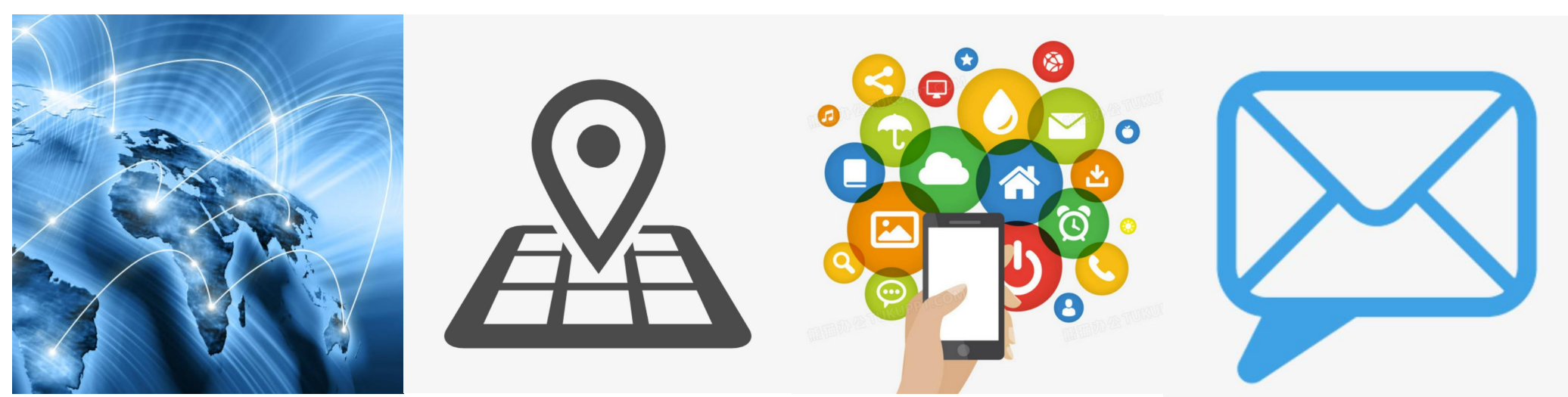
## Overview

### Background
- existing learned indexes fail to outperform traditional indexes when indexing **string keys**

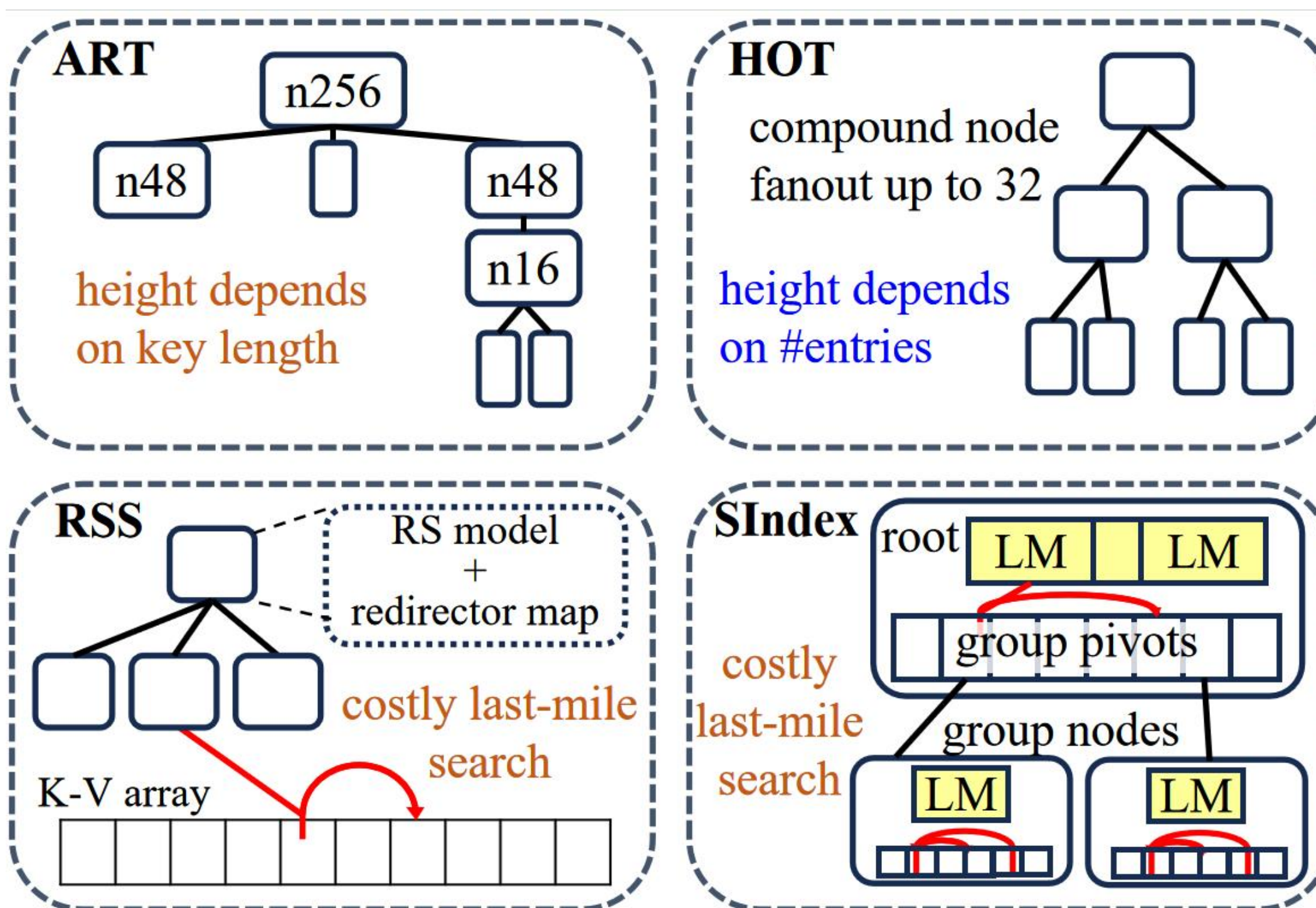### Characteristics of String keys
- long, variable-sized, with skewed prefixes

**string keys are common in real-world**

URL    address    user ID    email

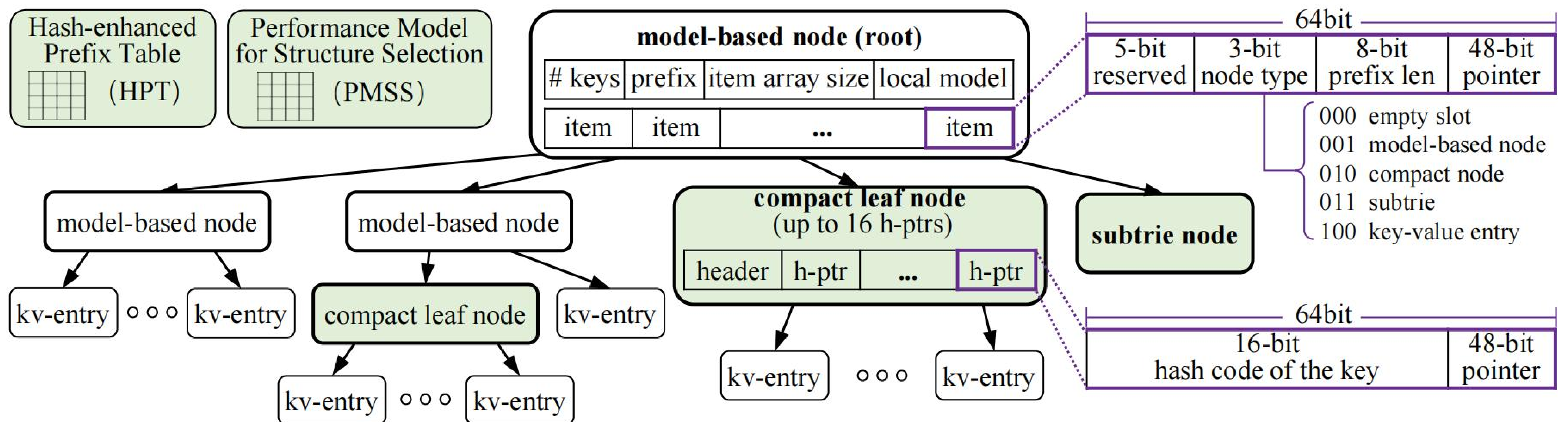### Existing Indexes Optimized for Strings

**ART**

height depends on key length

**HOT**

compound node fanout up to 32

height depends on #entries

**RSS**

RS model + redirector map

costly last-mile search

K-V array

**SIndex**

root — LM — LM

group pivots

group nodes

LM    LM

costly last-mile search

| | tree height | node search | last-mile search |
|---|---|---|---|
| ART | ✗ | ✓ | n/a |
| HOT | -- | ✓ | n/a |
| RSS | ✓ | ✓ | ✗ |
| SIndex | ✓ | ✓ | ✗ |
| LITS | ✓ | ✓ | n/a |

**our solution: LITS**
Learned Index with Hash-enhanced Prefix Table and Sub-tries

## LITS Structure

Hash-enhanced Prefix Table (HPT)

Performance Model for Structure Selection (PMSS)

**model-based node (root)**

| # keys | prefix | item array size | local model |

| item | item | ... | item |

model-based node

model-based node

**compact leaf node**
(up to 16 h-ptrs)

| header | h-ptr | ... | h-ptr |

**subtrie node**

kv-entry ○○○ kv-entry

compact leaf node    kv-entry

kv-entry ○○○ kv-entry

kv-entry ○○○ kv-entry

64bit

| 5-bit reserved | 3-bit node type | 8-bit prefix len | 48-bit pointer |

- 000 empty slot
- 001 model-based node
- 010 compact node
- 011 subtrie
- 100 key-value entry

64bit

| 16-bit hash code of the key | 48-bit pointer |

## Hash-enhanced Prefix Table

**CDF Model in LITS:** Hash-enhanced Prefix Table (HPT)

**Idea of HPT:** learn the pattern of a string data set by better approximating $prob(c|P)$

($prob(c|P)$: conditional probability of the next character being $c$ given the prefix $P$ )

target string: **bac**

| cdf\nprob | a | | b | | c | |
|---|---|---|---|---|---|---|
| 0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.7 | 0.3 |
| 1 | 0.2 | 0.1 | 0.3 | 0.5 | 0.8 | 0.2 |
| 2 | 0.1 | 0.3 | 0.4 | 0.1 | 0.5 | 0.5 |
| 3 | 0.4 | 0.4 | 0.8 | 0.1 | 0.8 | 0.1 |

b : 0 — 1

ba : 0.3 — 0.3+0.4

bac : 0.34 — 0.34+0.12
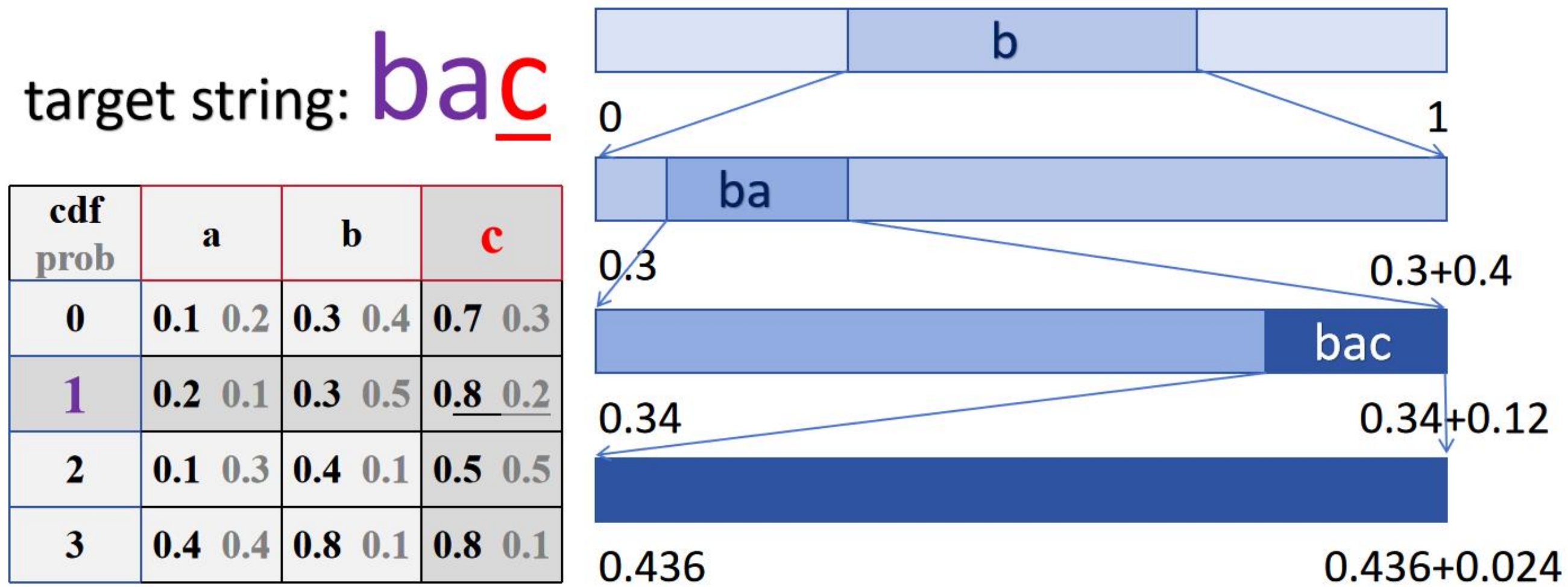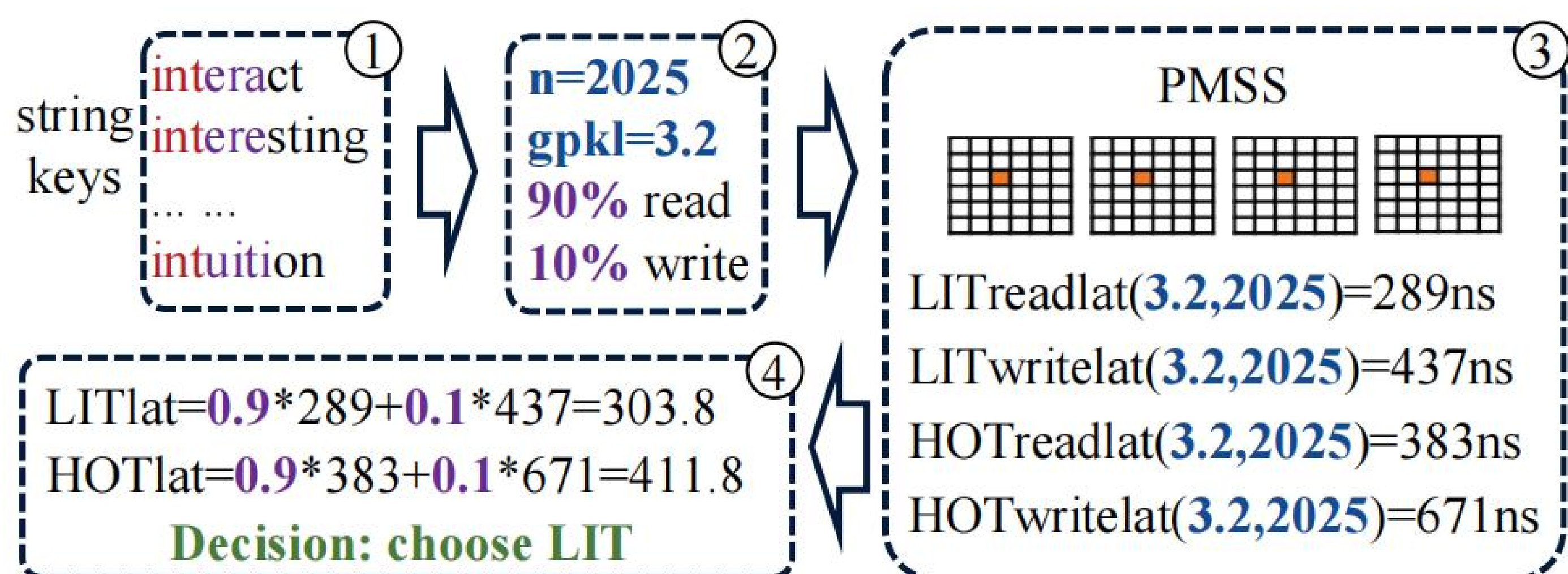
bac : 0.436 — 0.436+0.024

**Fig. An illustration of the CDF computation using the HPT for string "bac".**

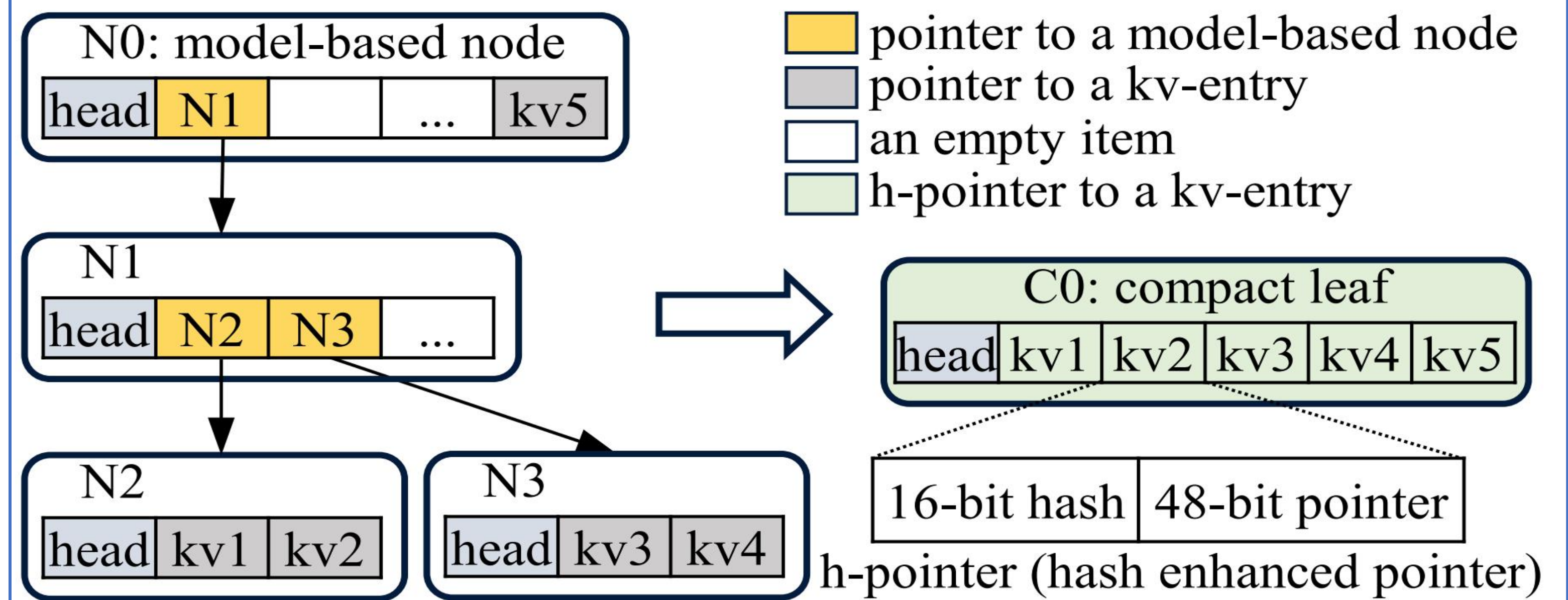## Performance Model for Structure Selection

**Hardness of String Data Set:** Group Partial Key Length (GPKL)
- GPKL reflects the hardness of modeling a string data set
- GPKL can be computed efficiently by reading the strings in one pass

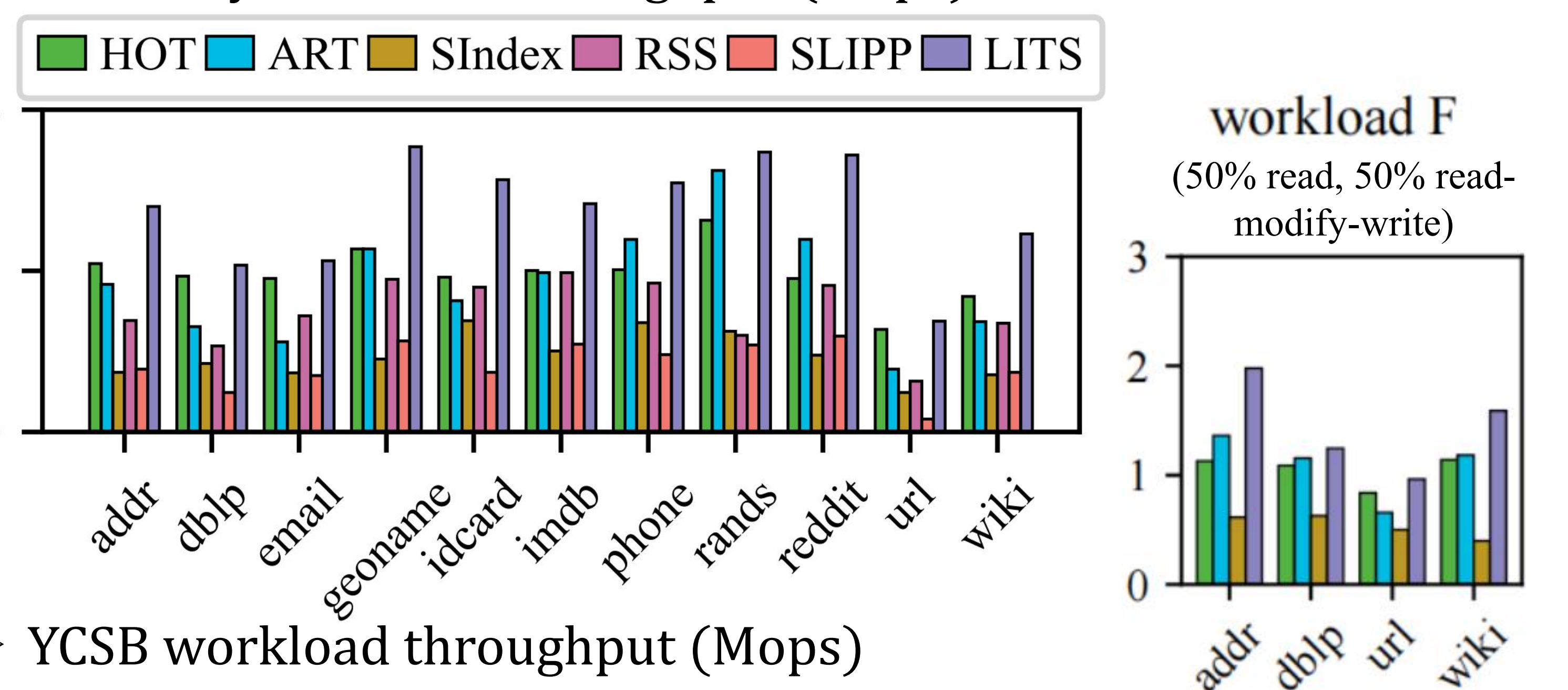**Idea of PMSS:** select the optimal structure for a sub-trie based on the hardness of a subset and offline benchmark tests
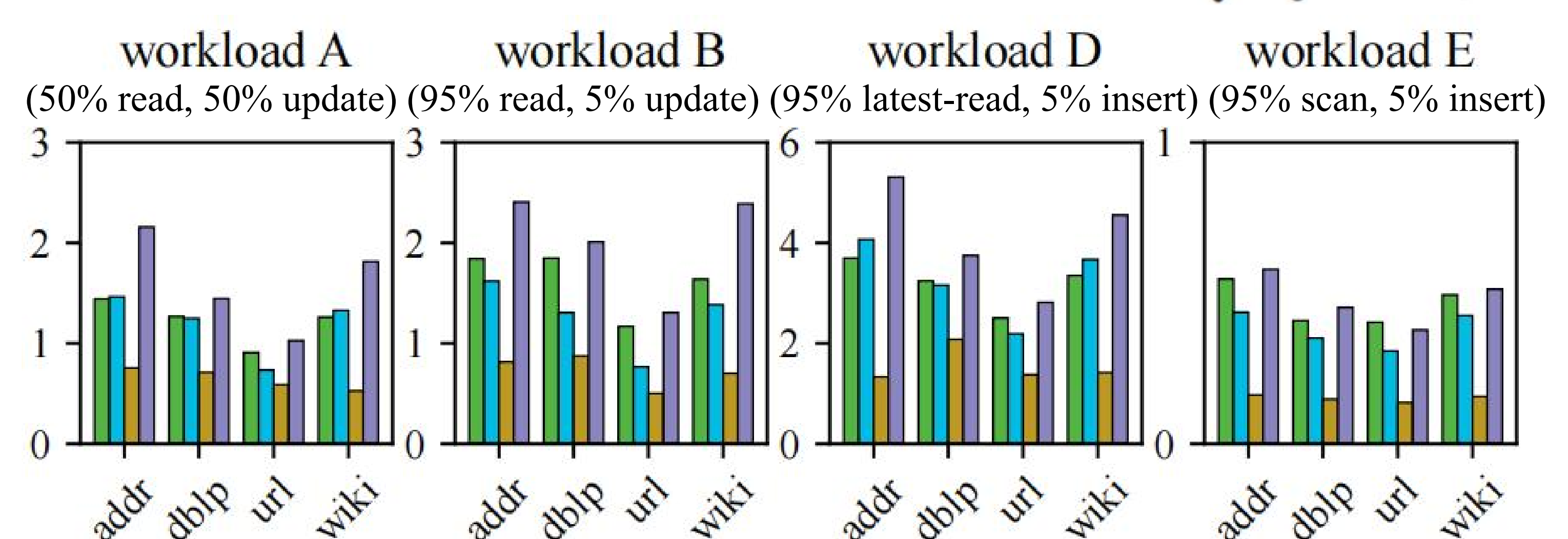
string keys: interact, interesting, ... ..., intuition ①

n=2025, gpkl=3.2, 90% read, 10% write ②

PMSS ③

LITreadlat(**3.2,2025**)=289ns
LITwritelat(**3.2,2025**)=437ns
HOTreadlat(**3.2,2025**)=383ns
HOTwritelat(**3.2,2025**)=671ns

LITlat=**0.9***289+**0.1***437=303.8
HOTlat=**0.9***383+**0.1***671=411.8
**Decision: choose LIT** ④

## Compact Leaf Node

N0: model-based node

| head | N1 | | ... | | kv5 |

N1

| head | N2 | N3 | | |

N2

| head | kv1 | kv2 |

N3

| head | kv3 | kv4 |

- 🟨 pointer to a model-based node
- ⬜ pointer to a kv-entry
- ⬜ an empty item
- 🟩 h-pointer to a kv-entry

C0: compact leaf

| head | kv1 | kv2 | kv3 | kv4 | kv5 |

| 16-bit hash | 48-bit pointer |

h-pointer (hash enhanced pointer)

## Evaluation

- ☐ 7 real-world string data sets and 4 synthetic string data sets
- ➤ read-only workload throughput (Mops)

HOT    ART    SIndex    RSS    SLIPP    LITS

(addr, dblp, email, geoname, idcard, imdb, phone, rands, reddit, url, wiki)

**workload F**
(50% read, 50% read-modify-write)

(addr, dblp, url, wiki)

- ➤ YCSB workload throughput (Mops)

**workload A**
(50% read, 50% update)

**workload B**
(95% read, 5% update)

**workload D**
(95% latest-read, 5% insert)

**workload E**
(95% scan, 5% insert)

(addr, dblp, url, wiki)

for more results please refer to our paper! 😊