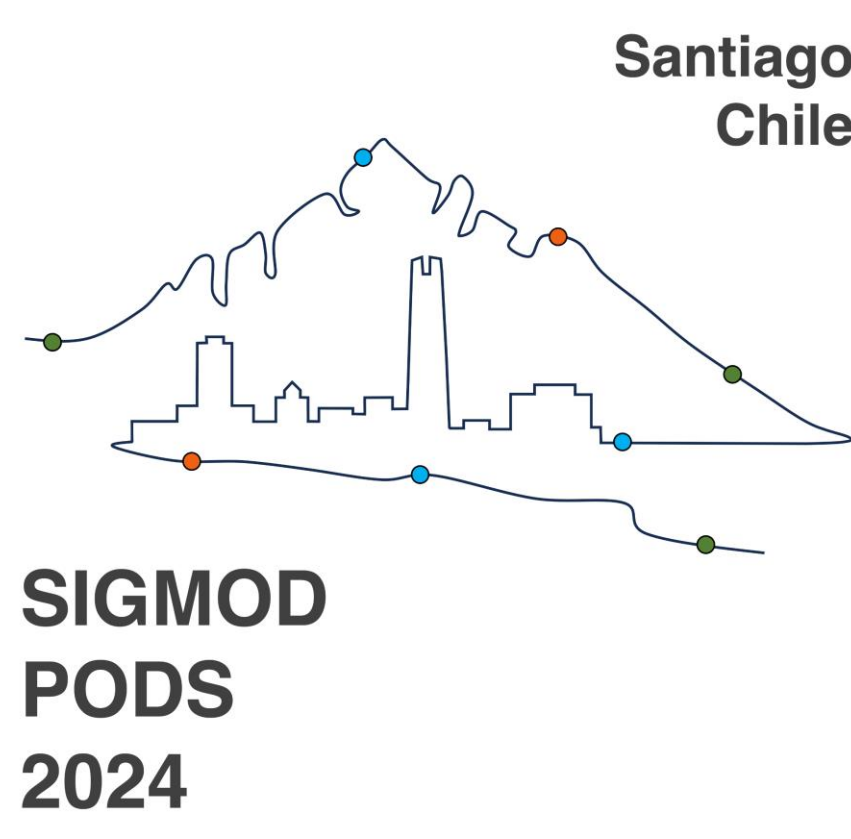# Can Learned Indexes be Built Efficiently?
# A Deep Dive into Sampling Trade-offs

**Minguk Choi**, Seehwan Yoo, and Jongmoo Choi
Dankook University, South Korea
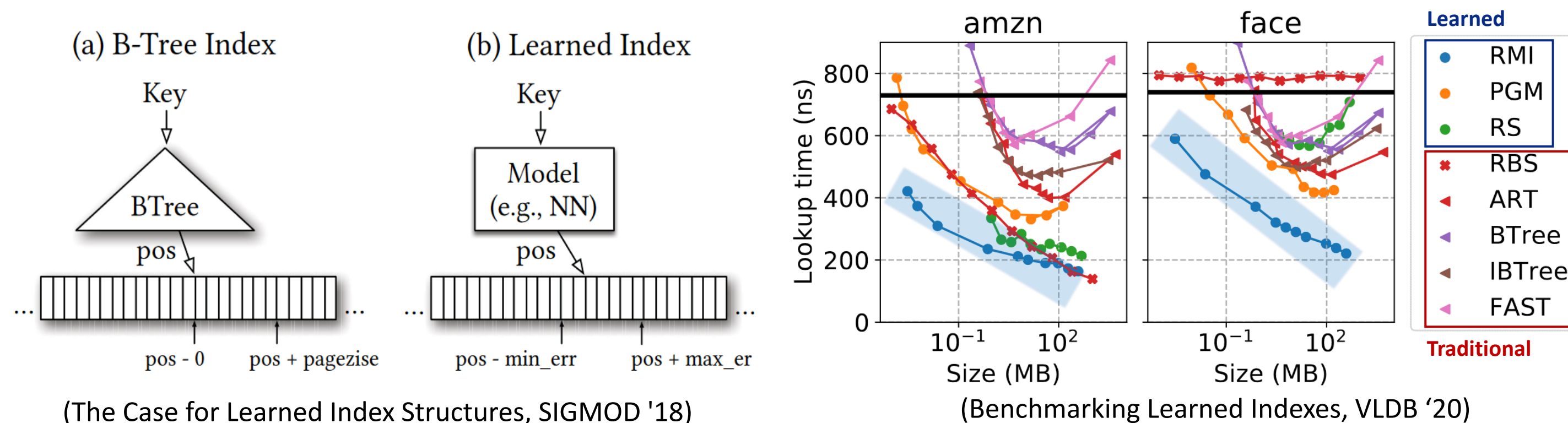{mgchoi, seehwan.yoo, jmchoi}@dankook.ac.kr

Seeking a Ph.D. Position for 2025

Santiago Chile
SIGMOD PODS 2024

## Background

- **Learned Index Structure**
  - Index structure employs machine learning techniques
  - View the index as a model that predicts the position of a key

- **Performance of Learned Index: Space-efficient**
  - Pareto optimal in terms of index size and lookup latency in read-only
    - No alternative exists that has both a smaller size and lower latency

(a) B-Tree Index
Key
BTree
pos
pos - 0    pos + pagesize

(b) Learned Index
Key
Model (e.g., NN)
pos
pos - min_err    pos + max_er

(The Case for Learned Index Structures, SIGMOD '18)

amzn / face
Learned: RMI, PGM, RS
Traditional: RBS, ART, BTree, IBTree, FAST
Lookup time (ns) vs Size (MB)
(Benchmarking Learned Indexes, VLDB '20)

## Motivation

- **Long Index Build Time**
  - Up to about 2,000x slower than traditional indexes
  - But still there are application where index build time is crucial (e.g., LSM-tree)

amzn
Lookup time (ns) vs Build time (ms)
Learned: RMI, PGM, RS
Traditional: CHT, RT, ART, BTree, IBTree
① For 650ns, up to 1,950x
② For 520ns, up to 187x
③ For 390ns, up to 8.3x

- **Why Building the Learned Index is Slow?**

  *Index build time = 1) Number of elements × 2) Per − element overhead*

  1) **Complete** traversal and training
  2) **Higher** per-element training overhead
     - Light-weight training model: RadixSpline (aiDM`20), Bourbon (OSDI`20)
       ➤ But it's **still longer** than traditional indexes

- **This study began with the question ...** 🤔
  Since the learned index uses the model,
  **Can't it learn efficiently even with less data?**

## Design

- **Our Approach: Sampling**

- **Challenges**
  1. Losing the error-bound property due to sampling loss
  2. Complex trade-offs in terms of model, index, and micro-architecture
  3. Absence of benchmark for sampling applied indexes

### 1. Error-bound Preserving Sample Learning Algorithm

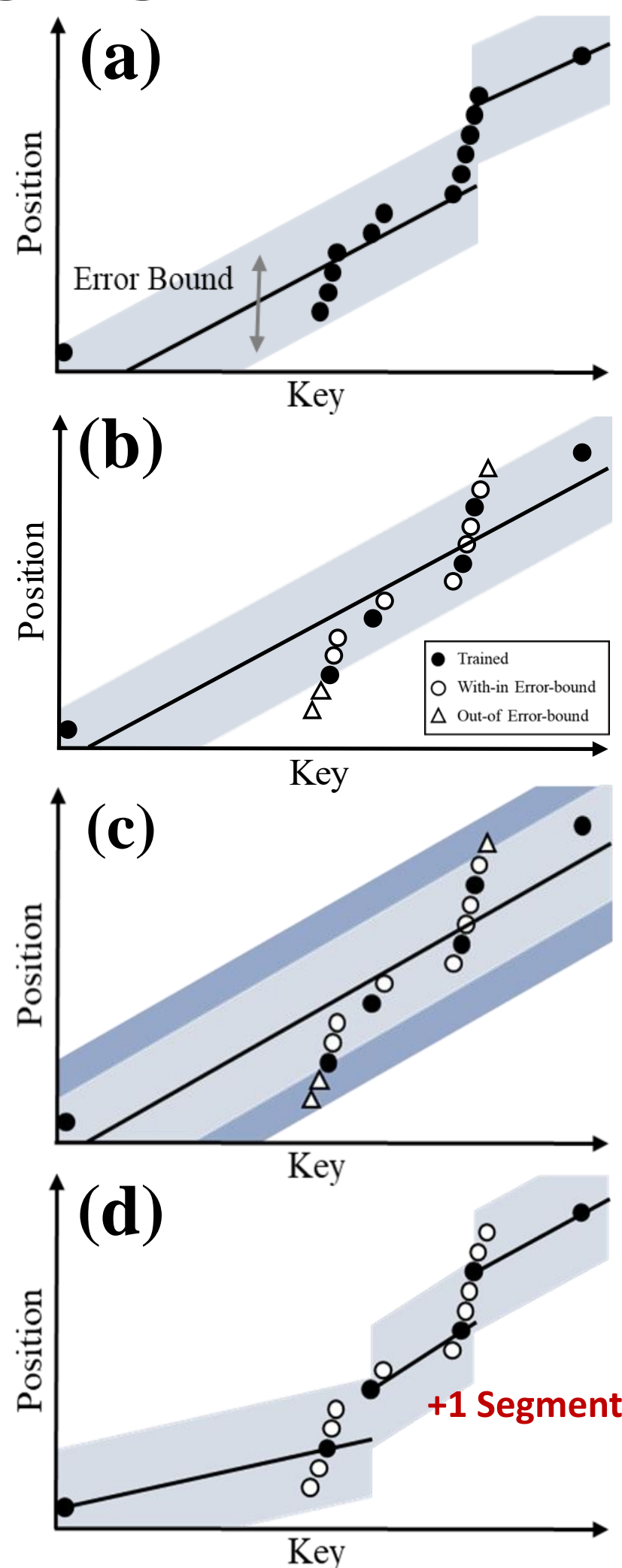- EB-PLA (Error-bounded Piece-wise Linear Approximation) Model
  (a) Train all keys with error-bound $\varepsilon$
    → $\forall k, Error(k) \leq \varepsilon$
  (b) Train sample $I^{th}$ keys with the error-bound $\varepsilon$
    ↛ $\forall k, Error(k) \leq \varepsilon$

- **Sample EB-PLA** Algorithm
  (c) Refine the error-bound due to sampling loss
    → $\forall k, Error(k) \leq \varepsilon'(= \varepsilon + I − 1)$
    ➤ Preserve the error-bound property
  (d) Replace the sample learning error-bound to
    $\delta (= \varepsilon − I + 1) → \forall k, Error(k) \leq \varepsilon$
    ➤ Preserve the error-bound ($\varepsilon$) by learning less data with smaller & stricter error-bound ($\delta$)

- Sample EB-Histogram

- PLR with Simple Linear Regression

(a) Position vs Key — Error Bound
(b) Position vs Key — Trained, With in Error-bound, Out-of Error-bound
(c) Position vs Key
(d) Position vs Key — +1 Segment

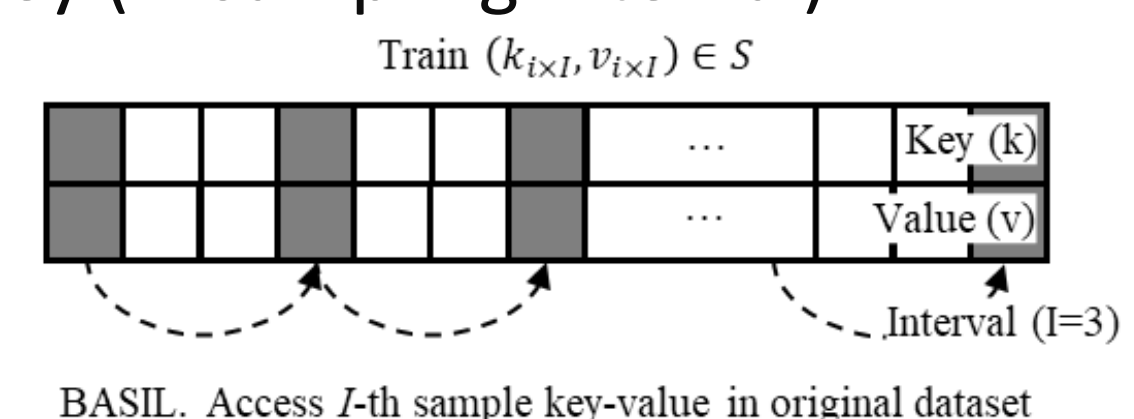## 2. Internal Changes due to Sampling

1) Dynamic Segmentation (Key range of each segment is different)
   - Aggressive sampling can increase the number of segments
2) Fixed Segmentation (Key range of each segment is equal)
   - Aggressive sampling can increase the number of under-fitting segments

(a) Simple Linear Regression (I = 1) — + Constant Segment
(b) Simple Linear Regression (I = 3)
(a) EB-Histogram (I = 1, $\varepsilon$ = 5) — Error Bound
(b) Sample EB-Histogram (I = 1, $\delta$ = 3, $\varepsilon$ = 5) — +4 Empty Bins, +1 Histogram — Trained, Not Trained

## 3. Unified Sampling Algorithm & Implementation

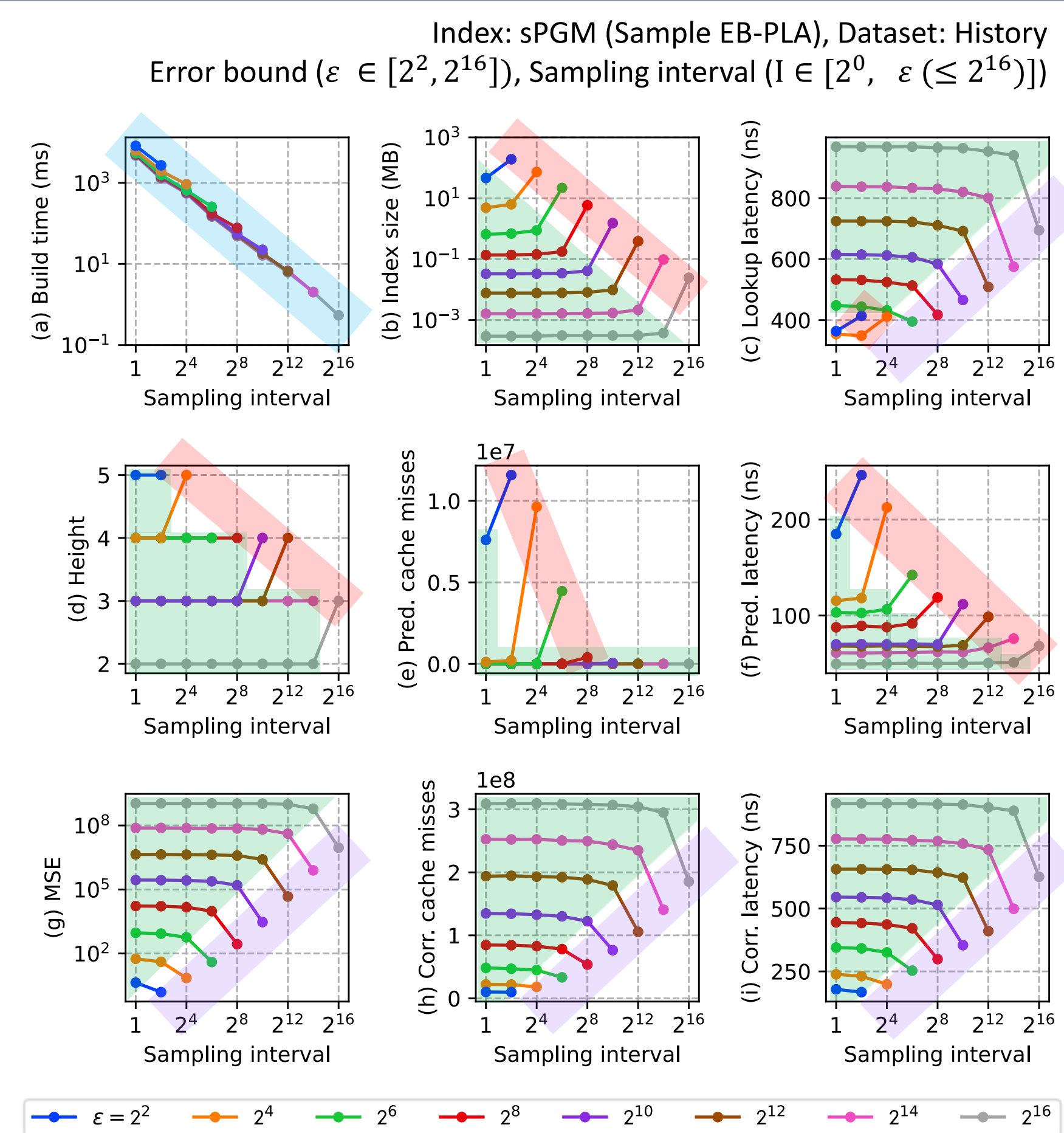🌱 **BASIL** (Benchmark of Sampling Applied Learned Indexes)
1) Unified Sampling Algorithm: Systematic Sampling
   - Extract every $I^{th}$ key form first to last key ($I$=sampling interval)
2) Unified Sampling Implementation
   - Index access and train only sample key-value data from entire dataset

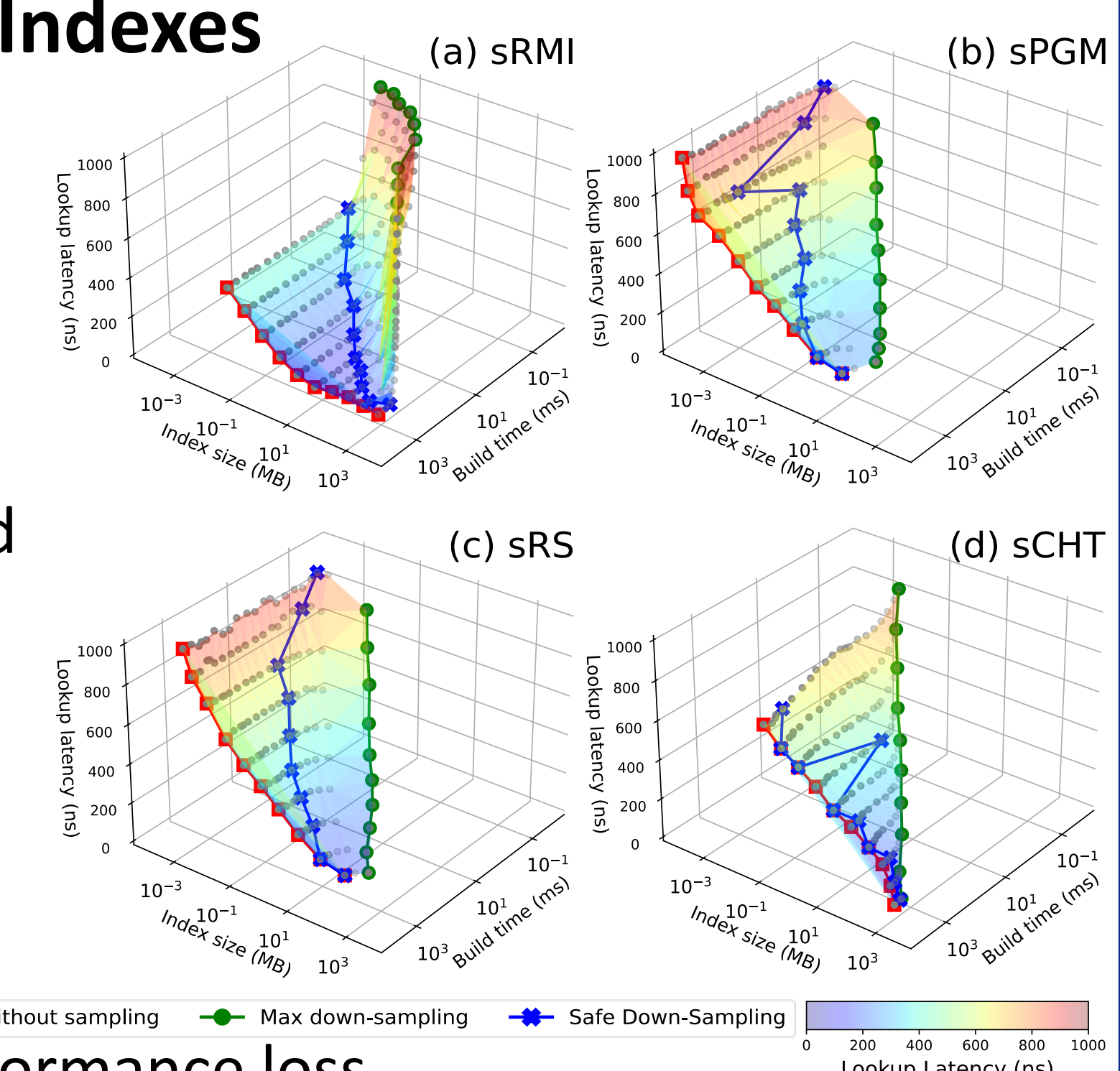Train $(k_{I \times I}, v_{I \times I}) \in S$
Key (k), Value (v), Interval (I=3)
BASIL. Access $I$-th sample key-value in original dataset

## Evaluation

### 1. Sampling Trade-offs

Index: sPGM (Sample EB-PLA), Dataset: History
Error bound ($\varepsilon \in [2^2, 2^{16}]$), Sampling interval ($I \in [2^0, \varepsilon (\leq 2^{16})]$)

- Sampling interval ($I$)↑ → (a) build time ↓
- Each error-bound ($\varepsilon$) has **threshold interval ($I^{TH}$)**
- Until $I^{TH}$, (b-i) rest of metrics remain consistent
- After $I^{TH}$, # of segments ↑
  → (b) Size ↑, (d) Height ↑
  → (e) Pred. cache miss ↑, (f) Pred. latency ↑
- After $I^{TH}$, # of segments ↑ → (g) MSE ↓ → (h) Corr. cache miss ↓, (i) Corr. latency ↓

(a) Build time (ms), (b) Index size (MB), (c) Lookup latency (ns) vs Sampling interval
(d) Height, (e) Pred. cache misses, (f) Pred. latency (ns) vs Sampling interval
(g) MSE, (h) Corr. cache misses, (i) Corr. latency (ns) vs Sampling interval
$\varepsilon = 2^2, 2^4, 2^6, 2^8, 2^{10}, 2^{12}, 2^{14}, 2^{16}$

### 2. Design Space of Learned Indexes

- Without sampling, absence of trade-offs between build, size, and lookup
- Sampling introduce trade-offs between build, size, and lookup
- **Broaden** design space of learned indexes from 2D to 3D

(a) sRMI, (b) sPGM, (c) sRS, (d) sCHT
Lookup latency (ns)

### 3. Build Speed-up

- Explore **Safe down-sampling**, where size & lookup latency increased by less than 5%
- Max build speedup without performance loss
  ➤ sRMI: 1/44,514, sPGM: 1/40,781, sRS: 1/14,479

Without sampling, Max down-sampling, Safe Down-Sampling

### 4. Pareto Analysis

- Can learned indexes be built more efficiently in terms of build time and lookup latency than traditional indexes through sampling?
- To the best of our knowledge, this is first to show that learned indexes are also **Pareto optimal** in terms of build time and (average and tail) lookup latency

history (Average), osm (Average), history (99.9%th), osm (99.9%th)
Avg. lookup latency (ns) / 99.9%th lookup latency (ns) vs Build time (ms)
Learned: sRMI, sPGM, sRS
Traditional: sCHT, sRT, sART, sBTree, sIBTree, BinarySearch