

Why are learned indexes so effective?



Paolo
Ferragina¹



Fabrizio
Lillo²



Giorgio
Vinciguerra¹

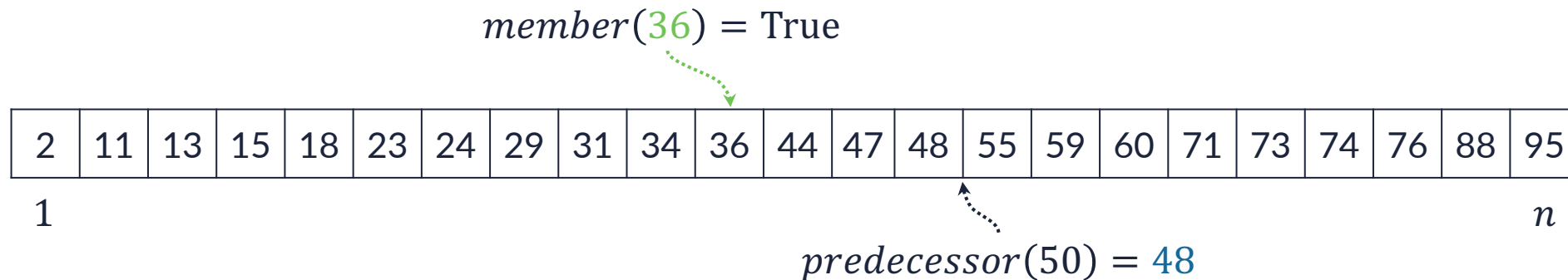


¹University of Pisa
²University of Bologna

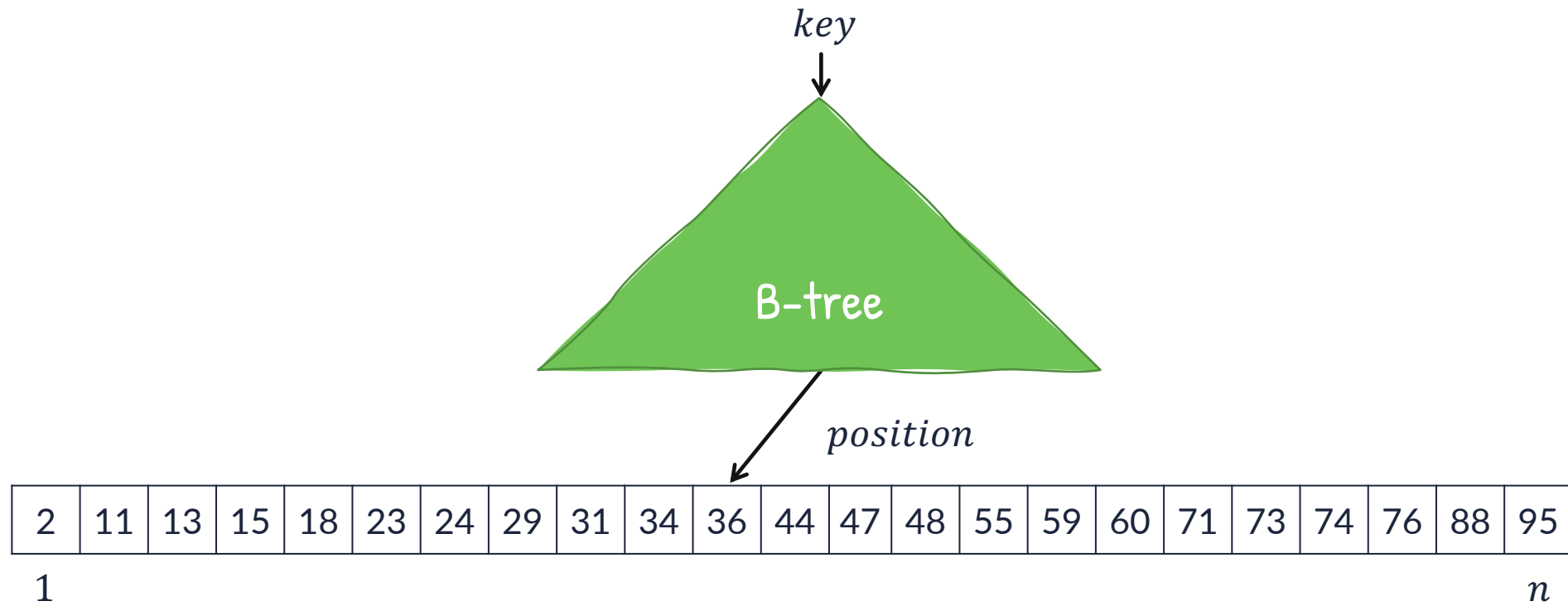


A classical problem in computer science

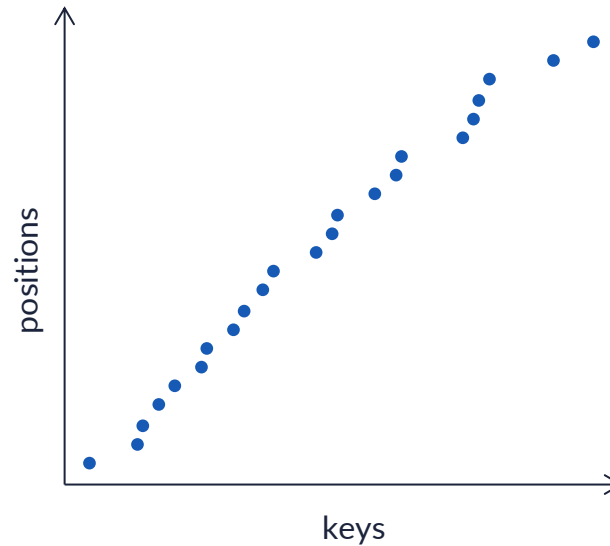
- Given a set of n sorted input keys (e.g. integers)
- Implement membership and predecessor queries
- Range queries in databases, conjunctive queries in search engines, IP lookup in routers...



Indexes

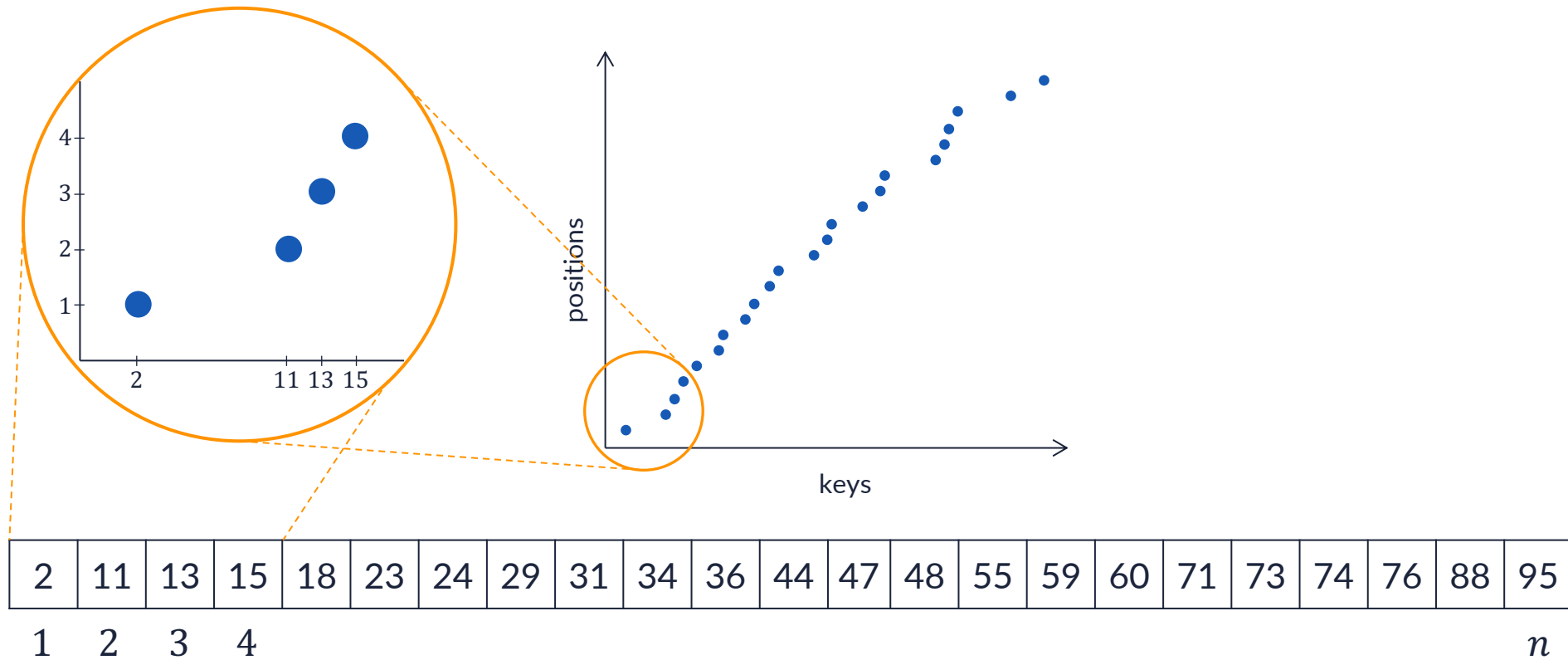


Input data as pairs (*key, position*)

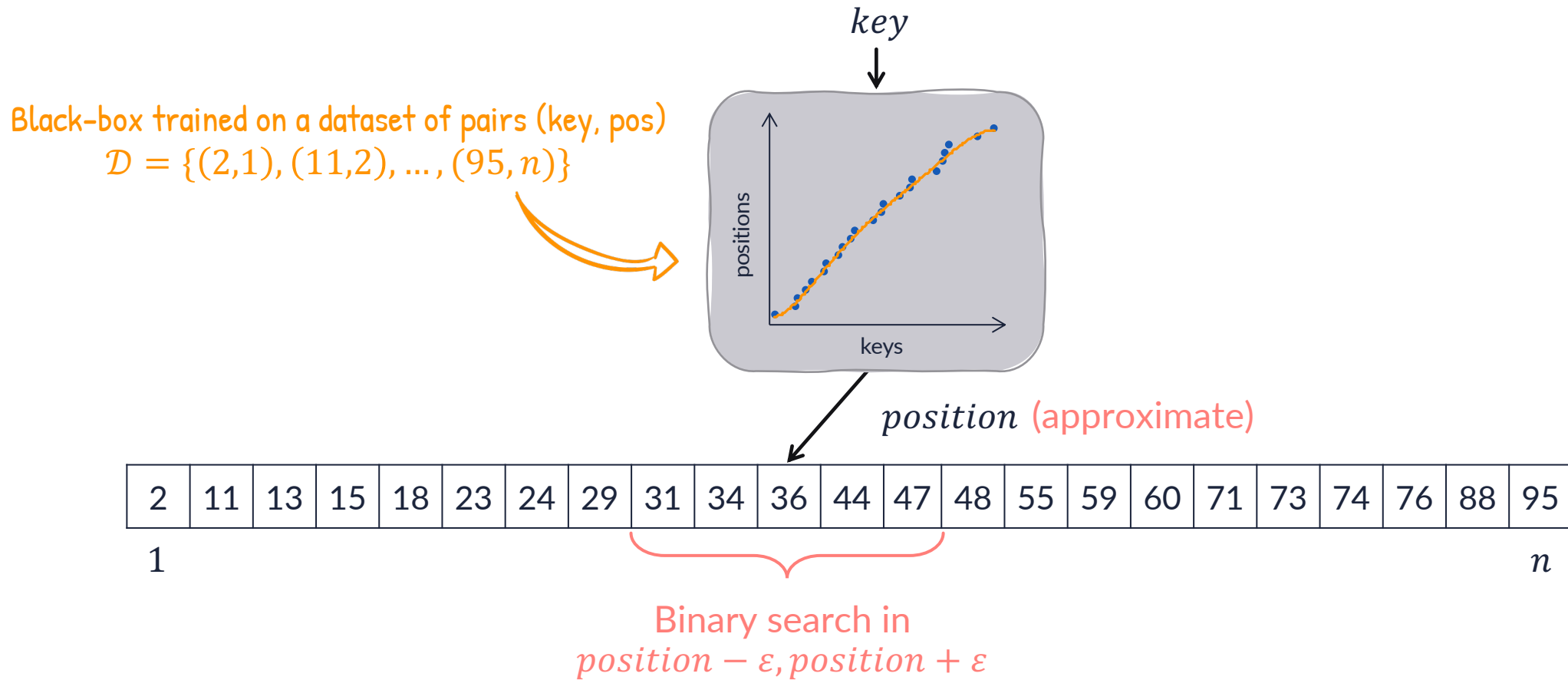


2	11	13	15	18	23	24	29	31	34	36	44	47	48	55	59	60	71	73	74	76	88	95	
1																							n

Input data as pairs (*key, position*)



Learned indexes



e.g. ϵ is of the order of 100–1000

The knowledge gap in learned indexes

Practice

Same query time of
traditional tree-based
indexes

vs
→

Theory

Same asymptotic query
time of traditional
tree-based indexes



Space improvements of
orders of magnitude,
from GBs to few MBs

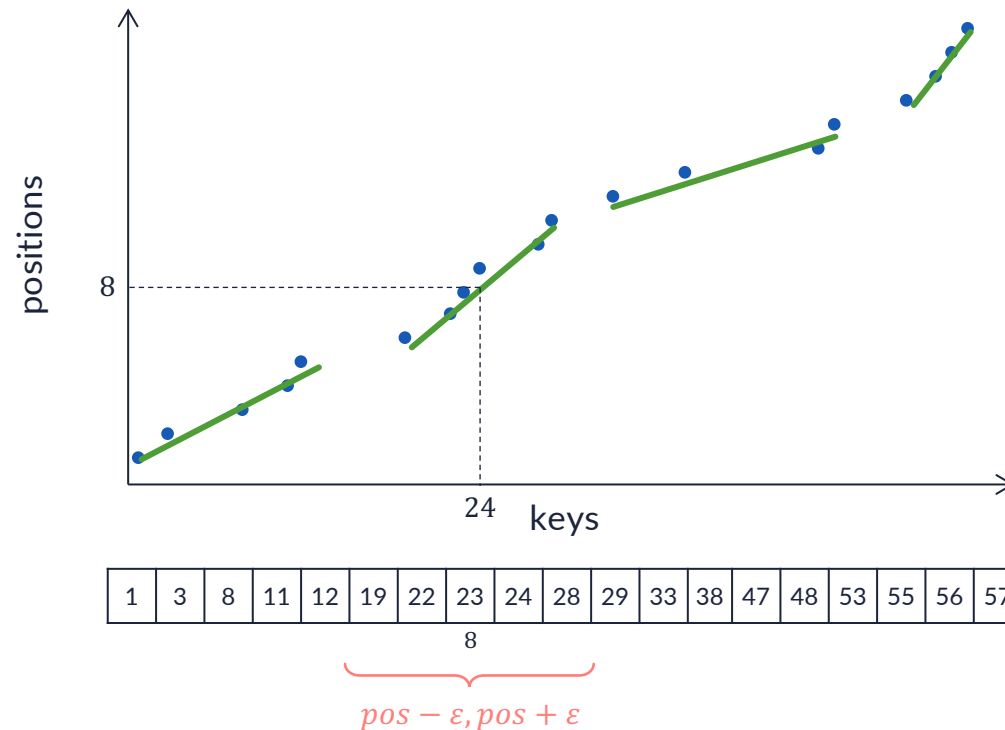
vs
→

Same asymptotic space
occupancy of traditional
tree-based indexes



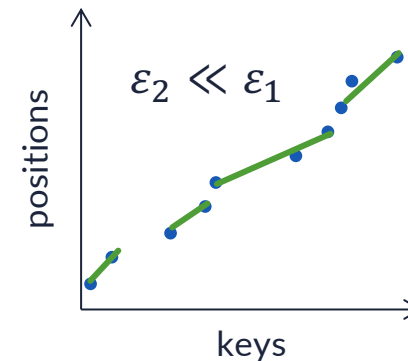
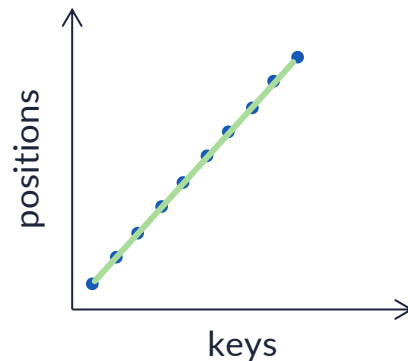
PGM-index: An optimal learned index

1. Fix a max error ε , e.g. so that keys in $[pos - \varepsilon, pos + \varepsilon]$ fit a cache-line
2. Find the smallest **Piecewise Linear ε -Approximation** (PLA)
3. Store triples (*first key, slope, intercept*) for each segment



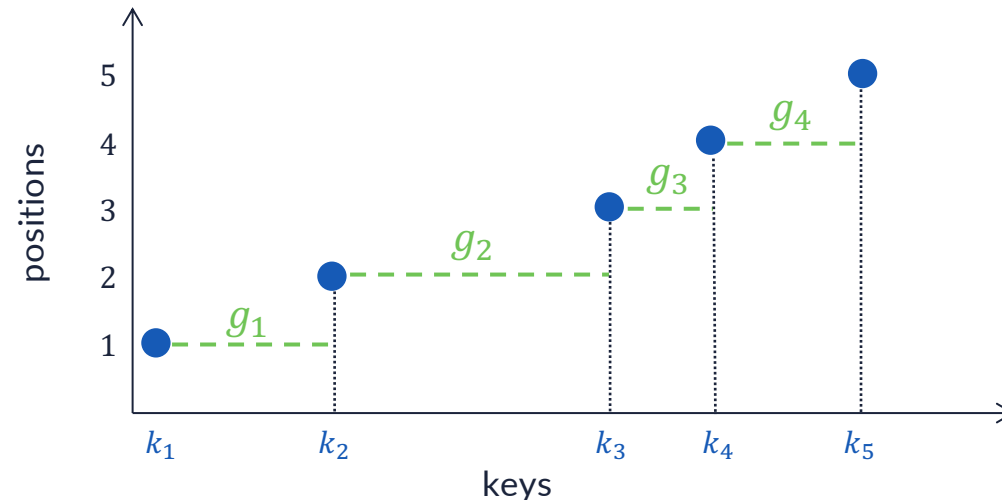
What is the space of learned indexes?

- Space occupancy \propto Number segments
- The number of segments depends on
 - The size of the input dataset
 - How the points (key, pos) map to the plane
 - The value ε , i.e. how much the approximation is precise



Model and assumptions

- Consider gaps $g_i = k_{i+1} - k_i$ between consecutive input keys
- Model the gaps as positive iid rvs that follow a distribution with finite mean μ and variance σ^2



The main result

Theorem. If ε is sufficiently larger than σ/μ , the expected number of keys covered by a segment with maximum error ε is

$$K = \frac{\mu^2}{\sigma^2} \varepsilon^2$$

and the number of segments on a dataset of size n is

$$\frac{n}{K}$$

with high probability.

The main consequence

The PGM-index achieves the same asymptotic query performance of a traditional ε -way tree-based index while improving its space from $\Theta(n/\varepsilon)$ to $O(n/\varepsilon^2)$

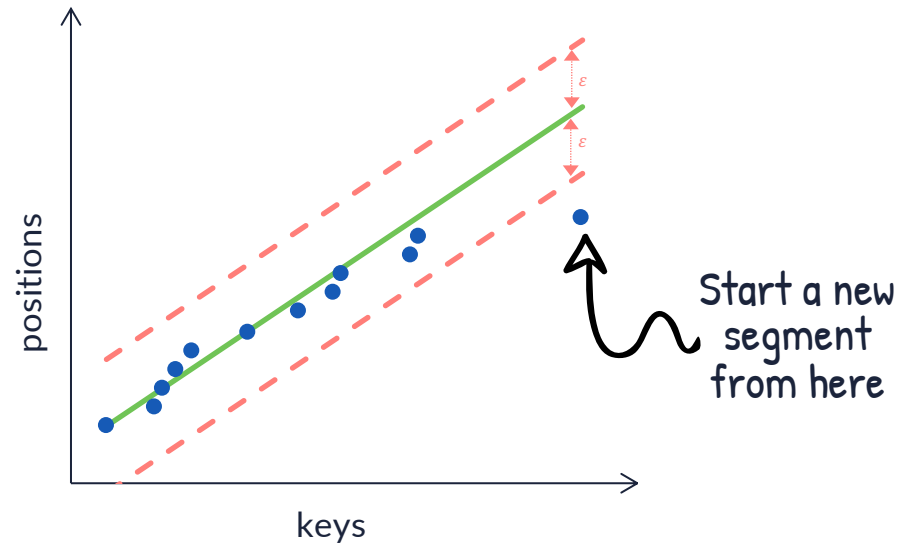


Learned indexes are provably
better than traditional indexes

(note that ε is of the order of 100-1000)

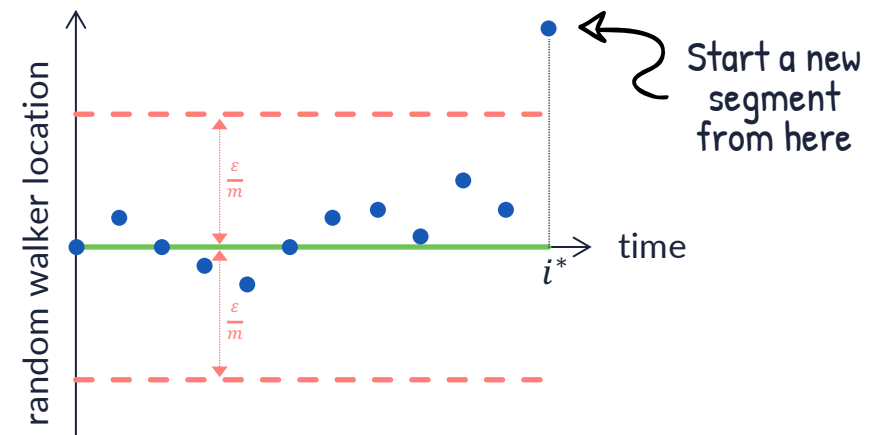
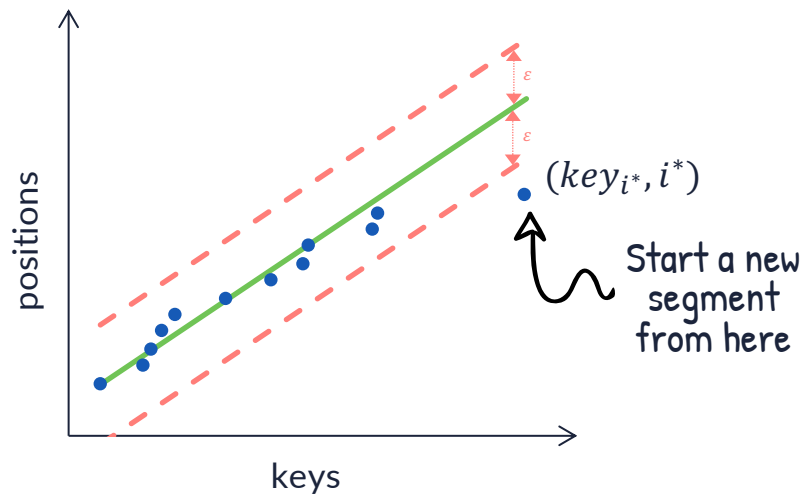
Sketch of the proof

1. Consider a **segment** on the stream of random gaps and the **two parallel lines** at distance ε
2. How many steps before a new segment is needed?



Sketch of the proof (2)

3. A discrete-time random walk, iid increments with mean μ
4. Compute the expectation of
$$i^* = \min\{i \in \mathbb{N} \mid (k_i, i) \text{ is outside the red strip}\}$$
i.e. the Mean Exit Time (MET) of the random walk
5. Show that the slope $m = 1/\mu$ maximises $E[i^*]$, giving $E[i^*] = (\mu^2/\sigma^2) \varepsilon^2$

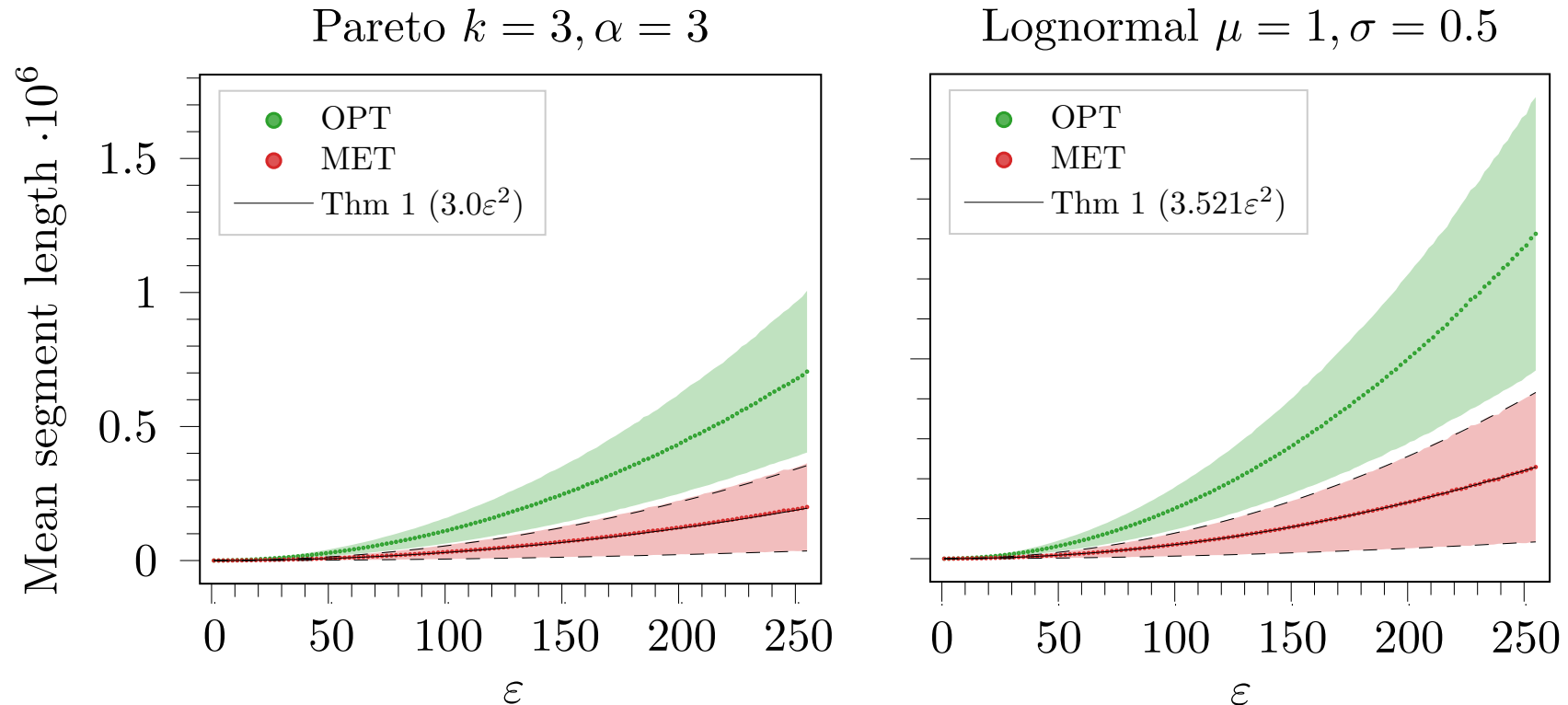


Simulations

1. Generate 10^7 random streams of gaps according to several probability distributions
2. Compute and average
 - I. The length of a segment found by the algorithm that computes the smallest PLA, adopted in the PGM-index
 - II. The exit time of the random walk

Simulations of $(\mu^2/\sigma^2)\varepsilon^2$

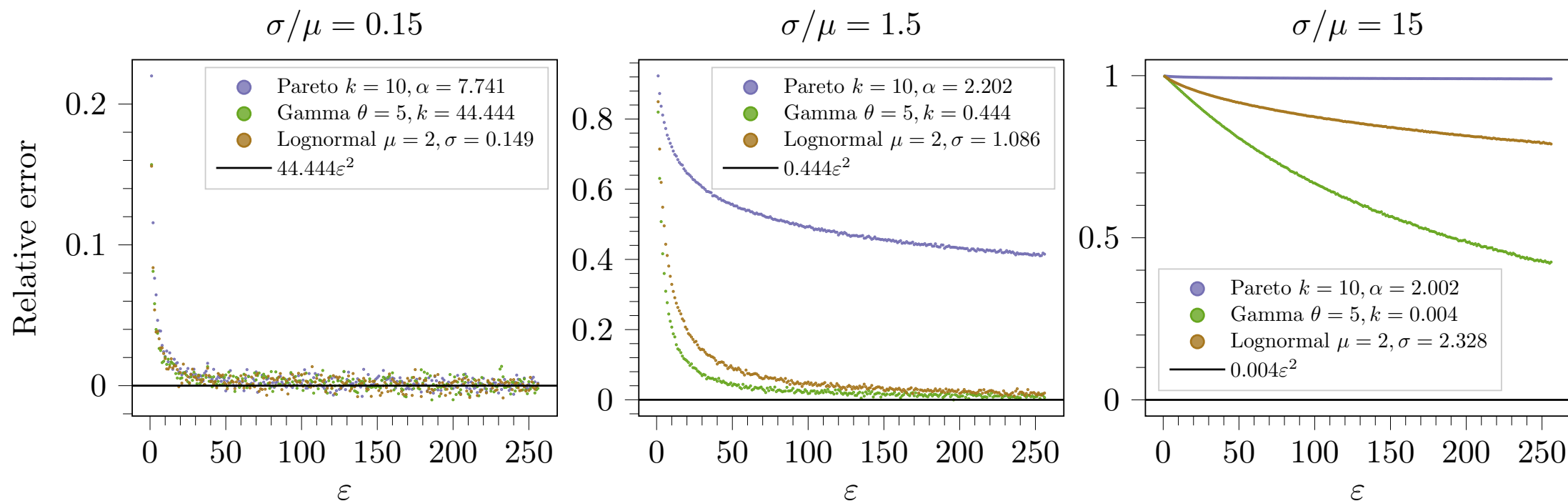
OPT = Average segment length in a PGM-index
MET = Mean exit time of the random walk



⇒ Both OPT and MET agree on the slope $1/\mu$, but OPT is more robust

More distributions in the paper

Stress test of “ ε sufficiently larger than σ/μ ”



Conclusions

- No theoretical grounds for the efficiency of learned indexes was known
- We have shown that on data with iid gaps, the mean segment length is $\Theta(\varepsilon^2)$
- The PGM-index takes $O(n/\varepsilon^2)$ space w.h.p., a quadratic improvement in ε over traditional indexes (ε is usually of the order of 100–1000)
- *Open problems:*
 1. Do the results still hold without the iid assumption on the gaps?
 2. Is the segment found by the optimal algorithm adopted in the PGM-index a constant factor longer than the one found by the random walker?