

# On the $k$ -Independence Required by Linear Probing and Minwise Independence

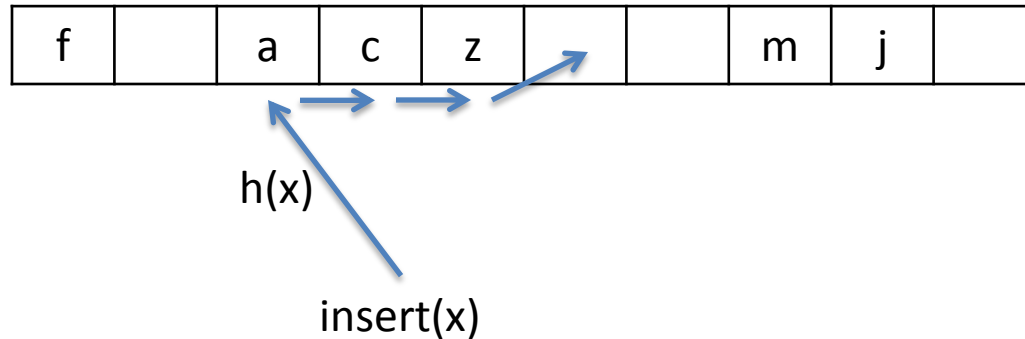
Mihai Pătrașcu

Mikkel Thorup



ICALP'10

# Linear Probing



[Knuth'63]  $E[\text{time of one operation}] = O(1)$   
“birth of algorithm analysis”

But assumes  $h$  is a truly random function  
 $\Rightarrow$  not an algorithm, but a *heuristic*

# Implementable Hash Functions

***k*-independence** [Wegman, Carter FOCS'79]

As we draw  $h$  from a family  $\mathcal{H}$ :

- uniformity:  $(\forall) x \in U, \quad h(x)$  uniform in  $[b]$
- independence:  $(\forall) x_1, \dots, x_k \in U, \quad h(x_1), \dots, h(x_k)$  i.i.d.

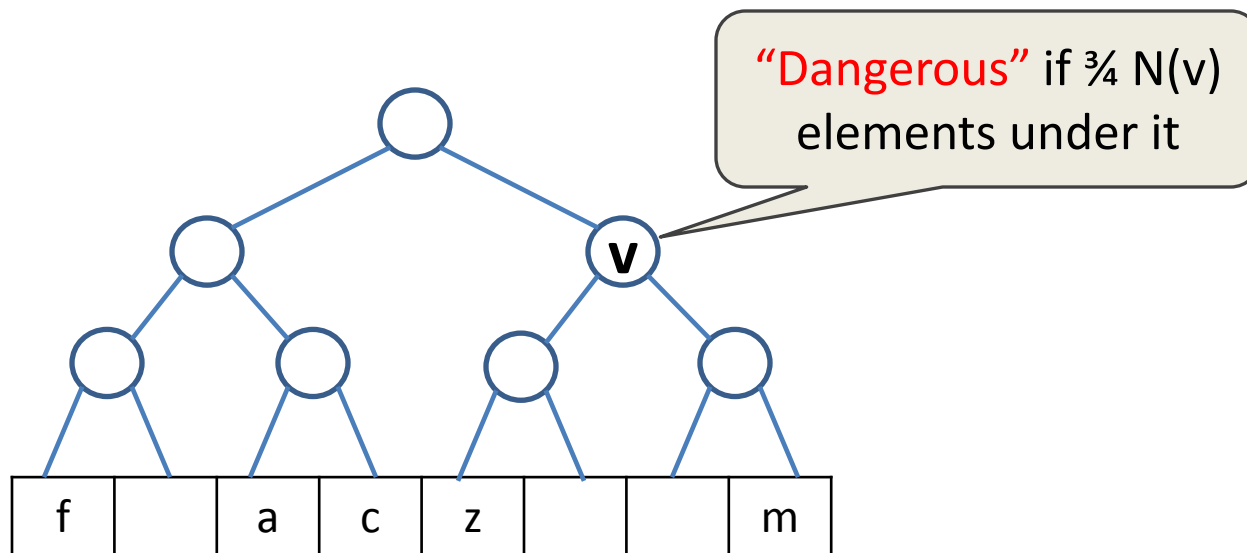
Possible implementation:

- let  $U =$  prime field
- draw  $a_0, \dots, a_{k-1} \in U$  randomly
- $h(x) = ( a_{k-1} x^{k-1} + \dots + a_1 x + a_0 ) \bmod b$

# Understanding Linear Probing

[Pagh, Pagh, Ružić STOC'07]

5-independence suffices!

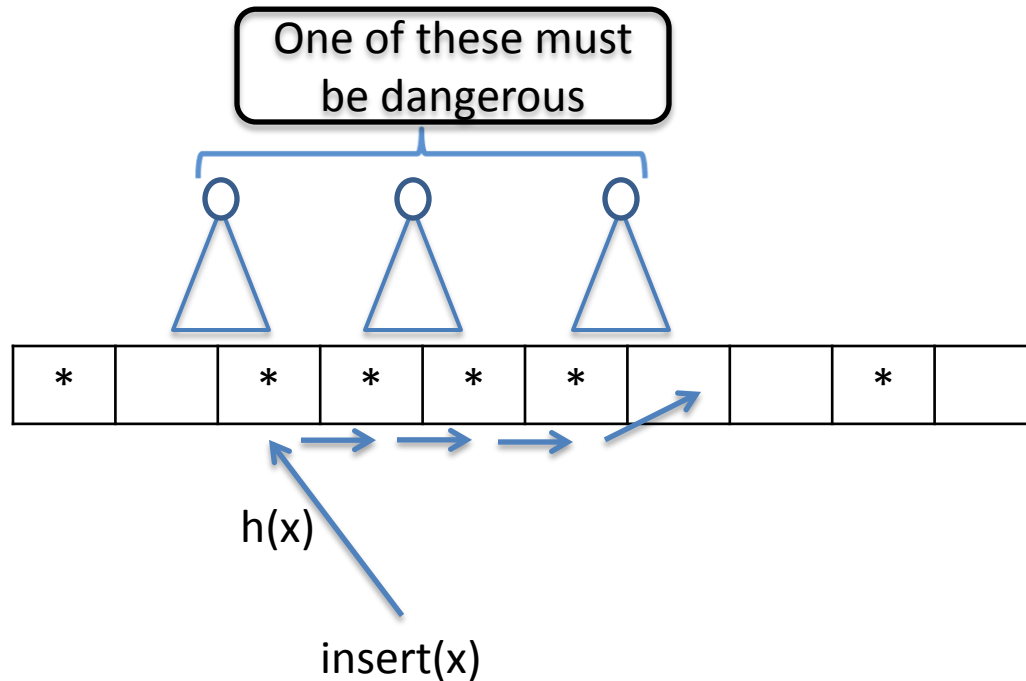


Let  $N(v) = \# \text{ leafs under } v = 2^{\text{level}(v)}$

# Understanding Linear Probing

*Main Lemma:* If  $h(x)$  is in a run of length  $2^k$

$\Rightarrow$  level  $k-1$  ancestor or a sibling must be dangerous



# Understanding Linear Probing

Look at “construction time” = time to insert  $n$  elements

Main lemma  $\Rightarrow$  construction time  $\leq \sum_{\text{dangerous}}$

Just a classic  
balls-in-bins analysis!

$E[\text{construction time}] \leq \sum_v [N(v)]^2 \cdot \Pr[v \text{ dangerous}]$

2-independence  $\Rightarrow$  Chebyshev bound

$\Rightarrow \Pr[v \text{ dangerous}] \leq 1/N(v)$

$\Rightarrow E[\text{construction time}] \leq \sum_v N(v) = O(n \lg n)$

# Understanding Linear Probing

Look at “construction time” = time to insert  $n$  elements

Main lemma  $\Rightarrow$  construction time  $\leq \sum_{\text{dangerous}}$

Just a classic  
balls-in-bins analysis!

$E[\text{construction time}] \leq \sum_v [N(v)]^2 \cdot \Pr[v \text{ dangerous}]$

4-independence  $\Rightarrow$  4<sup>th</sup> moment bound

$\Rightarrow \Pr[v \text{ dangerous}] \leq 1/[N(v)]^2$

$\Rightarrow E[\text{construction time}] \leq \sum_v O(1) = O(n)$

# Understanding Linear Probing

	$k=2$	$k=3$	$k=4$	$k=5$
Construction time	$O(n \lg n)$ $\Omega(n \lg n)$ [PPR]		$O(n)$	
Time/operation				



# Understanding Linear Probing

	$k=2$	$k=3$	$k=4$	$k=5$
Construction time	$O(n \lg n)$ $\Omega(n \lg n)$ [PPR]		$O(n)$	
Time/operation		$O(\lg n)$		$O(1)$

*One* query with  $k$ -independence

= keys arrange themselves by  $(k-1)$ -independence

+ the query hits a random location

# Understanding Linear Probing

	$k=2$	$k=3$	$k=4$	$k=5$
Construction time	$\Theta(n \lg n)$		$\Theta(n)$	
Time/operation	??	$O(\lg n)$		$\Theta(1)$

Do we really need “one more” for bounds / operation?

# Understanding Linear Probing

	$k=2$	$k=3$	$k=4$	$k=5$
Construction time	$\Theta(n \lg n)$		$\Theta(n)$	
Time/operation	$\Theta(\sqrt{n})$	$O(\lg n)$		$\Theta(1)$

Do we really need “one more” for bounds / operation?

YES.

Nasty 2-independent family such that:

- often,  $(\exists)$  run of  $\sqrt{n}$  elements;
- the query often falls in this bad run.

# Understanding Linear Probing

	$k=2$	$k=3$	$k=4$	$k=5$
Construction time	$\Theta(n \lg n)$	??	$\Theta(n)$	
Time/operation	$\Theta(\sqrt{n})$	$O(\lg n)$		$\Theta(1)$

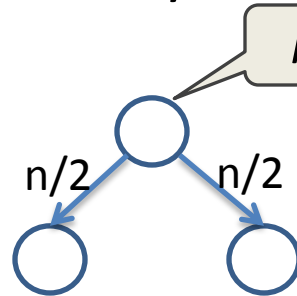
Could 3-independence help?

# Understanding Linear Probing

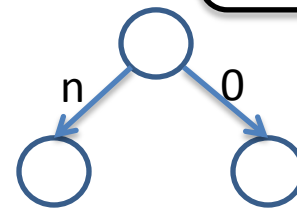
	$k=2$	$k=3$	$k=4$	$k=5$
Construction time	$\Theta(n \lg n)$	$\Omega(n \lg n)$	$\Theta(n)$	
Time/operation	$\Theta(\sqrt{n})$	$O(\lg n)$		$\Theta(1)$

Could 3-independence help?

Distribute keys down a tree:



Case 1



Case 2

If  $\Pr[\text{Case 2}] \approx 1/n$   
 $\Rightarrow$  3-independence!

cost  $\Omega(n^2)$

NO

# Understanding Linear Probing

	$k=2$	$k=3$	$k=4$	$k=5$
Construction time	$\Theta(n \lg n)$	$\Theta(n \lg n)$	$\Theta(n)$	
Time/operation	$\Theta(\sqrt{n})$	$\Theta(\lg n)$	??	$\Theta(1)$

Can both phenomena hit you simultaneously?

# Understanding Linear Probing

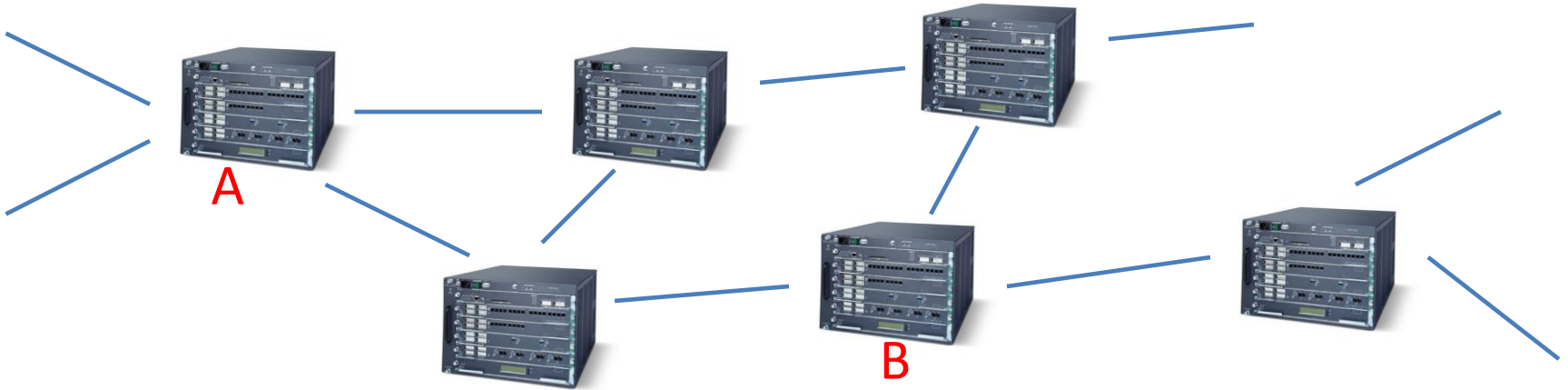
	$k=2$	$k=3$	$k=4$	$k=5$
Construction time	$\Theta(n \lg n)$	$\Theta(n \lg n)$	$\Theta(n)$	
Time/operation	$\Theta(\sqrt{n})$	$\Theta(\lg n)$	$O(\lg n)$ $\Omega(\lg n)$	$\Theta(1)$

Can both phenomena hit you simultaneously? YES

Apply the bad 3-independent distribution *only* on the query path  
 $\Rightarrow$  4-independent!

[Nasty proof ☹️]

# Minwise Independence



Problem: how many packets pass through both A and B?

Jaccard coefficient:  $|A \cap B| / |A \cup B|$

Algorithm:

- hash and keep  $\min h(A)$ ,  $\min h(B)$
- $\Pr[\min h(A) = \min h(B)] = |A \cap B| / |A \cup B|$
- repeat to estimate accurately



# Hashing Guarantees

Minwise Independence: for any  $S, x \in S$

$$\Pr[ h(x) = \min h(S) ] = 1 / |S|$$

Implies:  $\Pr[\min h(A) = \min h(B)] = |A \cap B| / |A \cup B|$

☹ Minwise independence not easy to obtain

# Hashing Guarantees

$\epsilon$ -Minwise Independence: for any  $S, x \in S$   
 $\Pr[ h(x) = \min h(S) ] = (1 \pm \epsilon) / |S|$

Implies:  $\Pr[\min h(A) = \min h(B)] = (1 \pm \epsilon) |A \cap B| / |A \cup B|$

[Indyk SODA'99] Any  $c \cdot \lg(1/\epsilon)$ -independent family  
is  $\epsilon$ -minwise independent

Here: Some  $c' \cdot \lg(1/\epsilon)$ -independent families  
are *not*  $\epsilon$ -minwise independent

# What it All Means

All our hash families are artificial

... we understand the *k-wise independence concept*

In practice:

- $(a * x) \gg \text{shift}$

More results:

- $\Omega(n \lg n)$ -construction for linear probing
- terrible minwise behavior

# What it All Means

All our hash families are artificial

... we understand the *k-wise independence concept*

In practice:

- $(a * x) \gg \text{shift}$
- tabulation-based hashing

Forthcoming paper:

Simple tabulation (3-wise independent) achieves

- linear probing in  $O(1)$  time (+ Chernoff concentration!)
- $o(1)$ -minwise independence

# What it All Means

All our hash families are artificial

... we understand the *k-wise independence concept*

In practice:

- $(a * x) \gg \text{shift}$
- tabulation-based hashing

The polynomial hash function:

- performance not understood
- but not too good in practice...

*The End*

Open problem: cuckoo hashing

6-independence needed [Cohen, Kane]

$O(\lg n)$ -independence suffices