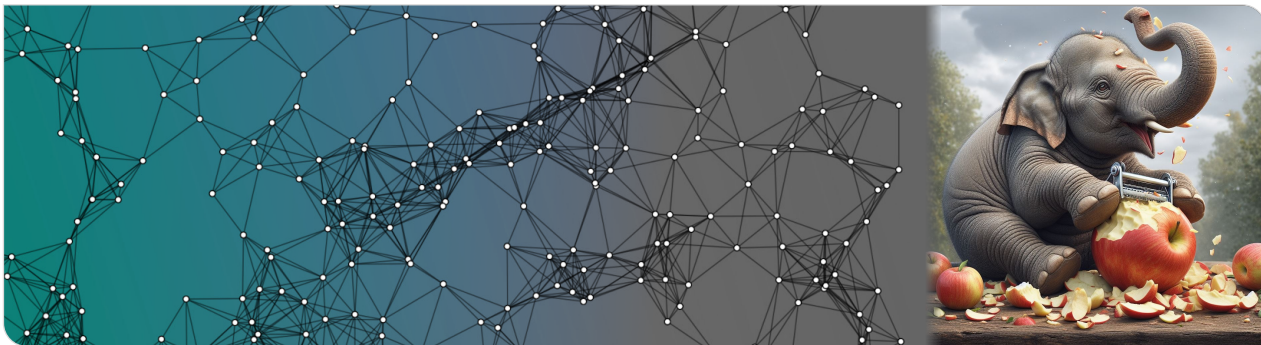


Probability and Computing – The Peeling Algorithm

Stefan Walzer, Maximilian Katzmann | WS 2023/2024



1. Cuckoo hashing with more than two hash functions

2. The Peeling Algorithm

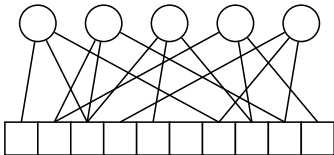
3. The Peeling Theorem

1. Cuckoo hashing with more than two hash functions

2. The Peeling Algorithm

3. The Peeling Theorem

Cuckoo Hashing with one table and k hash functions



$n \in \mathbb{N}$ keys
 $m \in \mathbb{N}$ table size
 $\alpha = \frac{n}{m}$ load factor
 $h_1, \dots, h_k \sim \mathcal{U}([m]^D)$ hash functions
 \hookrightarrow Could also use a separate table per hash function.

randomWalkInsert(x)

```

while  $x \neq \perp$  do // TODO: limit
  sample  $i \sim \mathcal{U}([k])$ 
  swap( $x, T[h_i(x)]$ )
  
```

(some improvements possible)

Theorem (without proof)

For each $k \in \mathbb{N}$ there is a **threshold** c_k^* such that:

- if $\alpha < c_k^*$ all keys can be placed with probability $1 - \mathcal{O}(\frac{1}{m})$.
- if $\alpha > c_k^*$ **not** all keys can be placed with probability $1 - \mathcal{O}(\frac{1}{m})$.

$$c_2^* = \frac{1}{2}, \quad c_3^* \approx 0.92, \quad c_4^* \approx 0.98, \dots$$

Conjecture

If $\alpha < c_k^*$ then the expected number of steps of successful insertions is $\mathcal{O}(1)$.

\hookrightarrow several proof attempts for random walk and other algorithms exist, with *partial* success

Static Hash Table

$\text{construct}(S)$: builds table T with key set S

$\text{lookup}(x)$: checks if x is in T or not

↔ no insertions or deletions after construction!

Constructing cuckoo hash tables:

- solved by Khosla 2013: “Balls into Bins Made Faster”
- matching algorithm resembling preflow push
- expected running time $\mathcal{O}(n)$, finds placement whenever one exists
- not in this lecture

Greedily constructing cuckoo hash tables

- Peeling algorithm: simple but sophisticated analysis
- interesting applications beyond hash tables (see “retrieval” in next lecture)

1. Cuckoo hashing with more than two hash functions

2. The Peeling Algorithm

3. The Peeling Theorem

The Peeling Algorithm

$\text{constructByPeeling}(S \subseteq D, h_1, h_2, h_3 \in [m]^D)$

$T \leftarrow [\perp, \dots, \perp]$ // empty table of size m

while $\exists i \in [m] : \exists$ *exactly one* $x \in S : i \in \{h_1(x), h_2(x), h_3(x)\}$ **do**

 // x is only unplaced key that may be placed in i

$T[i] \leftarrow x$

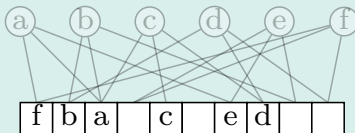
$S \leftarrow S \setminus \{x\}$

if $S = \emptyset$ **then**

return T

else

return NOT-PEELABLE



Exercise

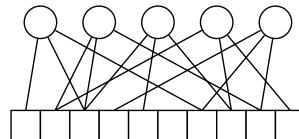
- Success of `constructByPeeling` does not depend on choices for i made by `while`.
- `constructByPeeling` can be implemented in linear time.

Cuckoo Graph and Peelability

- The **Cuckoo Graph** is the bipartite graph

$$G_{S, h_1, h_2, h_3} = (S, [m], \{(x, h_i(x)) \mid x \in S, i \in [3]\})$$

- Call G_{S, h_1, h_2, h_3} **peelable** if `constructByPeeling(S, h_1, h_2, h_3)` succeeds.
- If $h_1, h_2, h_3 \sim \mathcal{U}([m]^D)$ then the distribution of G_{S, h_1, h_2, h_3} does not depend on S . We then simply write $G_{m, \alpha m}$.
 - m \square -nodes and $\lfloor \alpha m \rfloor$ \circ -nodes
 - think: α is constant and $m \rightarrow \infty$.



Peeling simplified (not computing placement)

while \exists \square -node of degree 1 **do**
 \lfloor remove it and its incident \circ

G is peelable if and only if
this algorithm removes all \circ -nodes.

1. Cuckoo hashing with more than two hash functions

2. The Peeling Algorithm

3. The Peeling Theorem

Cuckoo hashing with more than two hash functions

○○○

The Peeling Algorithm

○○○

The Peeling Theorem

●○○○○○○○○○○○○○○○○○○○○

Peeling Theorem

Peeling Threshold

Let $c_3^\Delta = \min_{y \in [0,1]} \frac{y}{3(1-e^{-y})^2} \approx 0.81$.

Theorem (today's goal)

Let $\alpha < c_3^\Delta$. Then $\Pr[G_{m,\alpha m} \text{ is peelable}] = 1 - o(1)$.

Remark: More is known.

- For “ $\alpha < c_3^\Delta$ ” we get “peelable” with probability $1 - \mathcal{O}(1/m)$.
- For “ $\alpha > c_3^\Delta$ ” we get “not peelable” with probability $1 - \mathcal{O}(1/m)$.
- Corresponding thresholds c_k^Δ for $k \geq 3$ hash functions are also known.

Exercise: What about $k = 2$?

Peeling does not reliably work for $k = 2$ for any $\alpha > 0$.

Peeling Theorem: Proof outline

Theorem (today's goal)

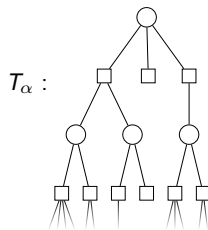
Let $\alpha < c_3^\Delta$. Then $\Pr[G_{m,\alpha m} \text{ is peelable}] = 1 - o(1)$.

Proof Idea

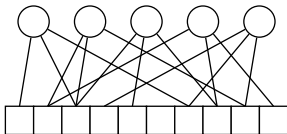
The random (possibly) infinite tree T_α can be peeled for $\alpha < c_3^\Delta$ and T_α is locally like $G_{m,\alpha m}$.

Steps

- I What is an infinite tree in general?
- II What is T_α in particular?
- III What does peeling mean in this setting?
- IV What role does c_3^Δ play?
- V What does it mean for T_α to be locally like $G_{m,\alpha m}$?
- VI What is the probability that a fixed key of $G_{m,\alpha m}$ is peeled?
- VII What is the probability that *all* keys of $G_{m,\alpha m}$ are peeled?



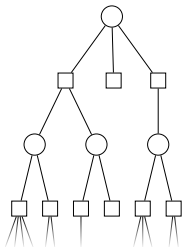
ii What is T_α in particular?



Observations for the finite Graph $G_{m, \alpha m}$

- each \bigcirc has 3 \square as neighbours (rare exception: $h_1(x), h_2(x), h_3(x)$ not distinct)
- each \square has random number X of \bigcirc as neighbours with $X \sim \text{Bin}(3n, \frac{1}{m}) = \text{Bin}(3\lfloor \alpha m \rfloor, \frac{1}{m})$. In an exercise you'll show

$$\Pr[X = i] \xrightarrow{m \rightarrow \infty} \Pr_{Y \sim \text{Pois}(3\alpha)} [Y = i].$$



Definition of the (possibly) infinite random tree T_α

- root is \bigcirc and has three \square as children
- each \square has random number of \bigcirc children, sampled $\text{Pois}(3\alpha)$ (independently for each \square).
- each non-root \bigcirc has two \square as children.

Remark: T_α is finite with positive probability > 0 , e.g. when the first three $\text{Pois}(3\alpha)$ random variables come out as 0. But T_α is also infinite with positive probability.

iii What does peeling mean in this setting? (2)

Observation

Let $q_R = \Pr[\text{root survives when peeling } T_\alpha^R]$.
The values q_R are decreasing in R .

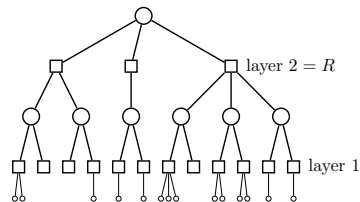
Peeling Algorithm

while \exists *childless* \square -*node* **do**
 \square remove it and its incident \circ

Proof.

Assume when peeling T_α^R the sequence $\vec{x} = (x_1, \dots, x_k)$ is a valid sequence of \square -node choices. Then \vec{x} is also valid when peeling T_α^{R+1} .

peeling T_α^R removes the root \Rightarrow peeling T_α^{R+1} removes the root
root survives when peeling $T_\alpha^{R+1} \Rightarrow$ peeling T_α^R removes the root
 $q_{R+1} \leq q_R$ □

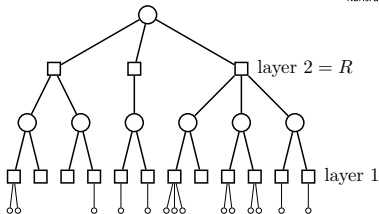


iii What does peeling mean in this setting? (3)

Peeling T_α^R bottom up

```

for  $i = 1$  to  $R$  do //  $\square$ -layers bottom to top
  for each  $\square$ -node  $v$  in layer  $i$  do
    if  $v$  has no children then
      remove  $v$  and its parent  $\circ$ 
  
```



Survival probabilities $p_i := \Pr[\square\text{-node in layer } i \text{ is not peeled}]$

$$\begin{aligned}
 p_1 &= \Pr[\square\text{-node has } \geq 1 \text{ child}] \\
 &= \Pr_{Y \sim \text{Pois}(3\alpha)}[Y > 0] = 1 - e^{-3\alpha}. \\
 p_i &= \Pr[\text{layer } i \text{ } \square\text{-node } v \text{ has } \geq 1 \text{ surviving child}] \\
 &= \Pr_{X \sim \text{Pois}(3\alpha p_{i-1}^2)}[X > 0] = 1 - e^{-3\alpha p_{i-1}^2}.
 \end{aligned}$$

\square -survival probabilities. With $p_0 := 1$ we have

$$p_i = \begin{cases} 1 & \text{if } i = 0 \\ 1 - e^{-3\alpha p_{i-1}^2} & \text{if } i = 1, 2, \dots \end{cases}$$

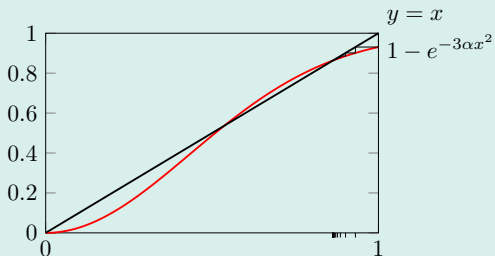
Moreover: $q_R := \Pr[\text{root survives}] = p_R^3$.

iv What role does $c_3^\Delta \approx 0.81$ play?

$$p_i = \begin{cases} 1 & \text{if } i = 0 \\ 1 - e^{-3\alpha p_{i-1}^2} & \text{if } i = 1, 2, \dots \end{cases}$$

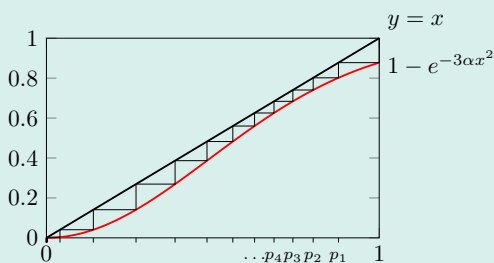
\hookrightarrow consider $f(x) = 1 - e^{-3\alpha x^2}$

Case 1: $\exists x > 0 : f(x) = x$.



$$\Rightarrow \lim_{i \rightarrow \infty} p_i = x^* = \max\{x \in [0, 1] \mid f(x) = x\}.$$

Case 2: $\forall x \in (0, 1] : f(x) < x$



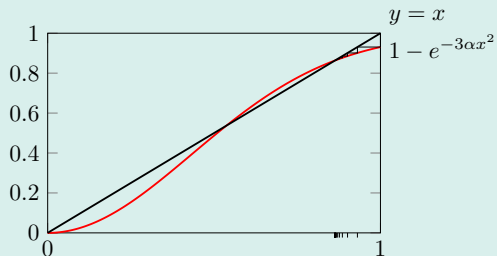
$$\Rightarrow \lim_{i \rightarrow \infty} p_i = 0.$$

iv What role does $c_3^\Delta \approx 0.81$ play?

$$p_i = \begin{cases} 1 & \text{if } i = 0 \\ 1 - e^{-3\alpha p_{i-1}^2} & \text{if } i = 1, 2, \dots \end{cases}$$

\hookrightarrow consider $f(x) = 1 - e^{-3\alpha x^2}$

Case 1: $\exists x > 0 : f(x) = x$.



$$\Rightarrow \lim_{i \rightarrow \infty} p_i = x^* = \max\{x \in [0, 1] \mid f(x) = x\}.$$

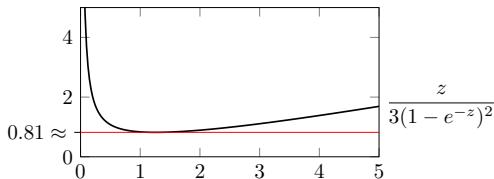
$$\text{Case 1} \Leftrightarrow \exists x > 0 : x = 1 - e^{-3\alpha x^2}$$

$$\Leftrightarrow \exists x > 0 : x^2 = (1 - e^{-3\alpha x^2})^2$$

$$\Leftrightarrow \exists z > 0 : \frac{z}{3\alpha} = (1 - e^{-z})^2 // z = 3\alpha x^2$$

$$\Leftrightarrow \exists z > 0 : \alpha = \frac{z}{3(1 - e^{-z})^2}$$

$$\Leftrightarrow \alpha \geq \min_{z > 0} \frac{z}{3(1 - e^{-z})^2} =: c_3^\Delta \approx 0.81$$



iv Interim Conclusion: What we learned about peeling T_α

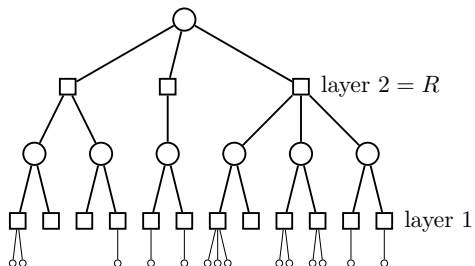
Lemma

For $\alpha < c_3^\Delta \approx 0.81$ we have

■ $\lim_{i \rightarrow \infty} p_i = 0.$

■ $\lim_{R \rightarrow \infty} q_R = \lim_{R \rightarrow \infty} p_R^3 = 0.$

“Root rarely survives for large R .”



v What does it mean for T_α to be locally like $G_{m,\alpha m}$?

Neighbourhoods in T_α and G

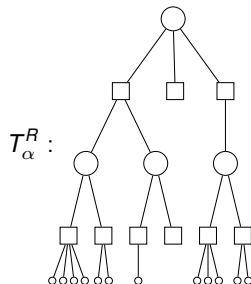
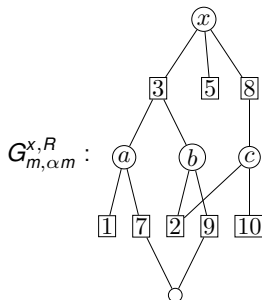
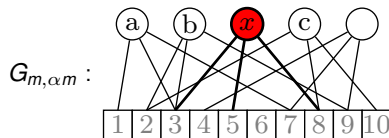
Let $R \in \mathbb{N}$. We consider

- T_α^R as before and
- for any fixed $x \in S$ the subgraph $G_{m,\alpha m}^{x,R}$ of $G_{m,\alpha m}$ induced by all nodes with distance at most $2R$ from x .

Lemma

For any $R \in \mathbb{N}$, the **distribution** of $G_{m,\alpha m}^{x,R}$ converges the distribution of T_α^R , i.e.

$$\forall T : \lim_{m \rightarrow \infty} \Pr[G_{m,\alpha m}^{x,R} = T] = \Pr[T_\alpha^R = T].$$

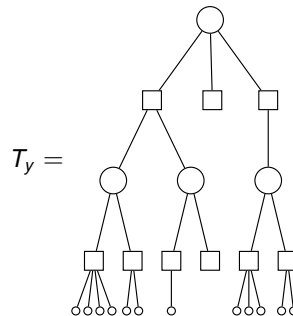


Lemma

Let T_y be a possible outcome of T_α^R given by a finite sequence $y = (y_1, \dots, y_k) \in \mathbb{N}_0^k$ specifying the number of children of \square -nodes in level order. Then

$$\Pr[T_\alpha^R = T_y] = \prod_{i=1}^k \Pr_{Y \sim \text{Pois}(3\alpha)} [Y = y_i].$$

e.g. for $y = (2, 0, 1, 4, 2, 1, 0, 3, 2)$:



v No cycles in $G_{m,\alpha m}^{x,R}$

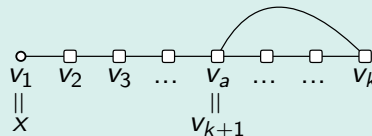
Lemma

Assume $R = \mathcal{O}(1)$. The probability that $G_{m,\alpha m}^{x,R}$ contains a cycle is $\mathcal{O}(1/m)$.

Proof.

If $G_{m,\alpha m}^{x,R}$ contains a cycle then we have

- a sequence $(v_1 = x, v_2, \dots, v_k, v_{k+1} = v_a)$ of nodes with $a \in [k]$
- of length $k \leq 4R$ (consider BFS tree for x and additional edge in it)
- for each $i \in \{1, \dots, k\}$ an index $j_i \in \{1, 2, 3\}$ of the hash function connecting v_i and v_{i+1} . (If $a = k - 1$ then $j_k \neq j_{k-1}$.)



$\Pr[\exists \text{ cycle in } G_{m,\alpha m}^{x,R}] \leq \Pr[\exists 2 \leq k \leq 4R : \exists v_2, \dots, v_k : \exists a \in [k] : \exists j_1, \dots, j_k \in [3] : \forall i \in [k] : h_{j_i} \text{ connects } v_i \text{ to } v_{i+1}]$

$$\leq \sum_{k=2}^{4R} \sum_{v_2, \dots, v_k} \sum_{a=1}^k \sum_{j_1, \dots, j_k} \prod_{i=1}^k \Pr[h_{j_i} \text{ connects } v_i \text{ to } v_{i+1}] \leq \sum_{k=2}^{4R} (\max\{m, n\})^{k-1} \cdot k \cdot 3^k \left(\frac{1}{m}\right)^k = \frac{1}{m} \sum_{k=2}^{4R} k \cdot 3^k = \mathcal{O}(1/m). \quad \square$$

v Distribution of $G_{m,\alpha m}^{x,R}$

Lemma

Let T_y be a possible outcome of T_α^R as before. Then

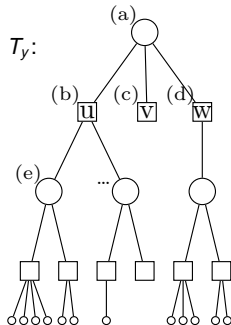
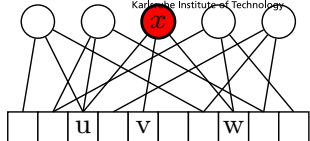
$$\Pr_{h_1, h_2, h_3 \sim \mathcal{U}([m]^D)} [G_{m,\alpha m}^{x,R} = T_y] \xrightarrow{m \rightarrow \infty} \prod_{i=1}^k \Pr_{Y \sim \text{Pois}(3\alpha)} [Y = y_i].$$

“Proof by example”, using T_y shown on the right.

The following things have to “go right” for $G_{m,\alpha m}^{x,R} = T_y$.

- a $h_1(x), h_2(x), h_3(x)$ pairwise distinct: probability $\xrightarrow{m \rightarrow \infty} 1$
 \hookrightarrow non-distinct would give cycle of length 2. Unlikely by lemma.

Note: $3 \lfloor \alpha m \rfloor - 3$ remaining hash values $\sim \mathcal{U}([m])$.



v Distribution of $G_{m,\alpha m}^{x,R}$

Lemma

Let T_y be a possible outcome of T_α^R as before. Then

$$\Pr_{h_1, h_2, h_3 \sim \mathcal{U}([m]^D)} [G_{m,\alpha m}^{x,R} = T_y] \xrightarrow{m \rightarrow \infty} \prod_{i=1}^k \Pr_{Y \sim \text{Pois}(3\alpha)} [Y = y_i].$$

“Proof by example”, using T_y shown on the right.

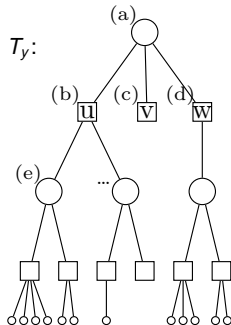
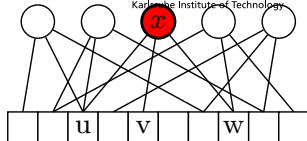
b Exactly $y_1 = 2$ of the remaining hash values are u .

$$\hookrightarrow \Pr_{Y \sim \text{Bin}(3 \lfloor \alpha m \rfloor - 3, \frac{1}{m})} [Y = 2] \xrightarrow{m \rightarrow \infty} \Pr_{Y \sim \text{Pois}(3\alpha)} [Y = 2]. \rightarrow \text{exercise}$$

Moreover: The two hash values must belong to 2 distinct keys. Probability $\xrightarrow{m \rightarrow \infty} 1$.

\hookrightarrow non-distinct would give cycle of length 2.

Note: The $3 \lfloor \alpha m \rfloor - 5$ remaining hash values are $\sim \mathcal{U}([m] \setminus \{u\})$. \rightarrow exercise



v Distribution of $G_{m,\alpha m}^{x,R}$

Lemma

Let T_y be a possible outcome of T_α^R as before. Then

$$\Pr_{h_1, h_2, h_3 \sim \mathcal{U}([m]^D)} [G_{m,\alpha m}^{x,R} = T_y] \xrightarrow{m \rightarrow \infty} \prod_{i=1}^k \Pr_{Y \sim \text{Pois}(3\alpha)} [Y = y_i].$$

“Proof by example”, using T_y shown on the right.

c None of the remaining hash values are v .

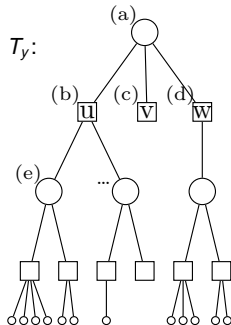
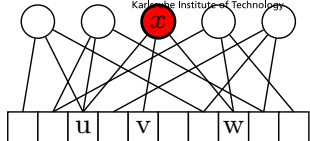
$$\hookrightarrow \Pr_{Y \sim \text{Bin}(3\lfloor \alpha m \rfloor - 5, \frac{1}{m-1})} [Y = 0] \xrightarrow{m \rightarrow \infty} \Pr_{Y \sim \text{Pois}(3\alpha)} [Y = 0].$$

Note: The $3\lfloor \alpha m \rfloor - 5$ remaining hash values are $\sim \mathcal{U}([m] \setminus \{u, v\})$.

d One of the remaining hash values is w .

$$\hookrightarrow \Pr_{Y \sim \text{Bin}(3\lfloor \alpha m \rfloor - 5, \frac{1}{m-2})} [Y = 1] \xrightarrow{m \rightarrow \infty} \Pr_{Y \sim \text{Pois}(3\alpha)} [Y = 1].$$

...



v Distribution of $G_{m,\alpha m}^{x,R}$

Lemma

Let T_y be a possible outcome of T_α^R as before. Then

$$\Pr_{h_1, h_2, h_3 \sim \mathcal{U}([m]^D)} [G_{m,\alpha m}^{x,R} = T_y] \xrightarrow{m \rightarrow \infty} \prod_{i=1}^k \Pr_{Y \sim \text{Pois}(3\alpha)} [Y = y_i].$$

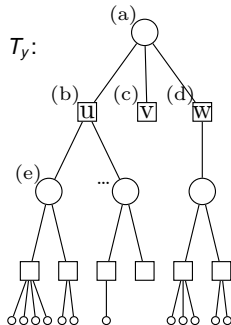
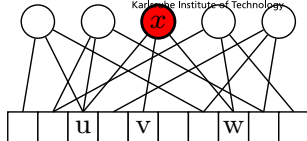
Proof sketch in general (some details omitted)

- General case at i -th \square -node. Want: probability that $G_{m,\alpha m}^{x,R}$ continues to match T_y . Note: T_y is fixed, so i and the number c_i of previously revealed hash values is bounded.

$$\Pr_{Y \sim \text{Bin}(3\lfloor \alpha m \rfloor - c_i, \frac{1}{m-i+1})} [Y = y_i] \xrightarrow{m \rightarrow \infty} \Pr_{Y \sim \text{Pois}(3\alpha)} [Y = y_i].$$

Moreover, those y_i hash values must belong to distinct fresh keys. Probability $\xrightarrow{m \rightarrow \infty} 1$
 \hookrightarrow otherwise we'd have a cycle.

- General case for \circ -node. The two children must be fresh: probability $\xrightarrow{m \rightarrow \infty} 1$
 \hookrightarrow otherwise there would be a cycle.



vi Probability that a specific key survives peeling

Lemma

Let $\alpha < c_3^\Delta$. Let x be any \bigcirc -node in $G_{m,\alpha m}$ as before (chosen before sampling the hash functions). Let

$$\mu_m := \Pr_{h_1, h_2, h_3 \sim \mathcal{U}([m]^D)} [x \text{ is removed when peeling } G_{m,\alpha m}].$$

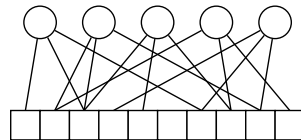
Then $\lim_{m \rightarrow \infty} \mu_m = 1$.

vii Proof of the Peeling Theorem

Theorem

Let $\alpha < c_3^\Delta$. Then

$$\Pr[G_{m,\alpha m} \text{ is peelable}] = 1 - o(1).$$



Proof

Let $n = \lfloor \alpha m \rfloor$ and $0 \leq s \leq n$ the number of \bigcirc nodes surviving peeling.

last lemma: each \bigcirc survives with probability $o(1)$.

linearity of expectation $\mathbb{E}[s] = n \cdot o(1) = o(n)$.

Exercise: $\Pr[s \in \{1, \dots, \delta n\}] = \mathcal{O}(1/m)$ if $\delta > 0$ is a small enough constant.

Markov: $\Pr[s > \delta n] \leq \frac{\mathbb{E}[s]}{\delta n} = \frac{o(n)}{\delta n} = o(1)$.

finally: $\Pr[s > 0] = \Pr[s \in \{1, \dots, \delta n\}] + \Pr[s > \delta n] = \mathcal{O}(1/m) + o(1) = o(1)$. \square

Conclusion

Peeling Process

- greedy algorithm for placing keys in cuckoo table
- works up to a load factor of $c_3^\Delta \approx 0.81$

We saw glimpses of important techniques

- *Local interactions in large graphs*. Also used in statistical physics.
- *Galton-Watson Processes / Trees*. Random processes related to T_α .
- *Local weak convergence*. How the finite graph $G_{m,\alpha m}$ is locally like T_α .

But wait, there's more!

- Further applications of peeling
 - retrieval data structures (next lecture)
 - perfect hash functions (next lecture)
 - set sketches
 - linear error correcting codes

- Cuckoo Hashing und der Schälalgorithmus
 - (Wie) kann man Cuckoo Hashing mit mehr als 2 Hashfunktionen aufziehen?
 - Welcher Vorteil ergibt sich im Vergleich zu 2 Hashfunktionen?
 - Wie funktioniert der Schälalgorithmus zur Platzierung von Schlüsseln in einer Cuckoo Hashtabelle?
 - Schälen lässt sich als einfacher Prozess auf Graphen auffassen. Wie?
 - Was besagt das Hauptresultat, das wir zum Schälprozess bewiesen haben?
- Beweis des Schälsatzes. *Mir ist klar, dass der Beweis äußerst kompliziert ist.*
 - Im Beweis haben zwei Graphen eine Rolle gespielt ein endlicher und ein (potentiell) unendlicher. Wie waren diese Graphen definiert?
 - Welcher Zusammenhang besteht zwischen der Verteilung der Knotengrade in T_α und $G_{m,\alpha m}$?