

# Practical Evaluation of LZ78 and LZW Tries

**Johannes Fischer**

*Dominik Köppl*

# LZ78 with LZ trie

senescence

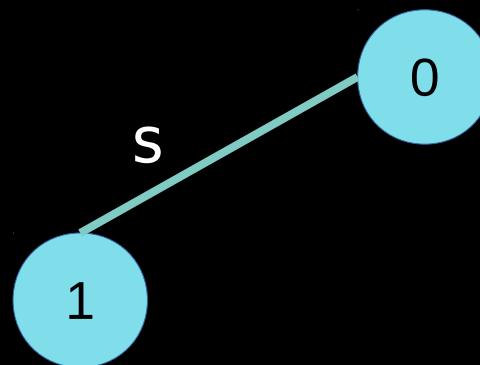
# LZ78 with LZ trie

0

senescence

# LZ78 with LZ trie

Senescence

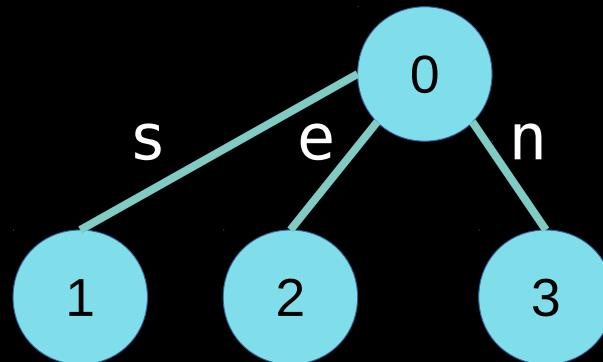


# LZ78 with LZ trie

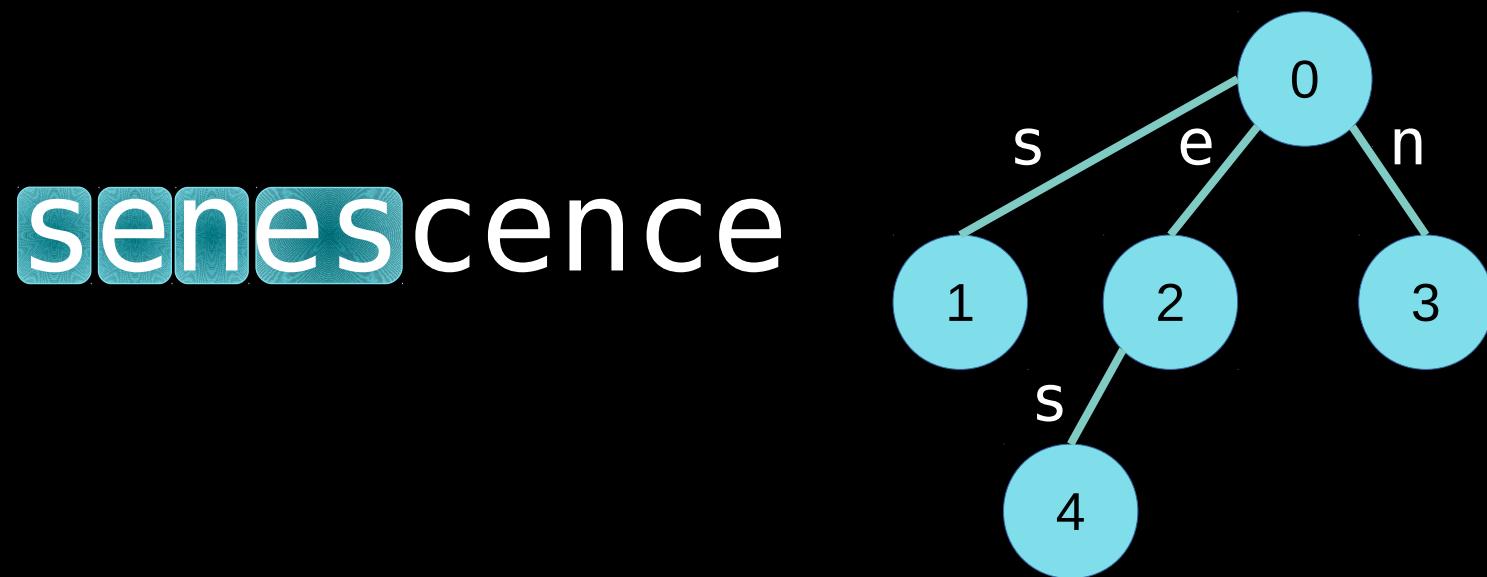


# LZ78 with LZ trie

senescence

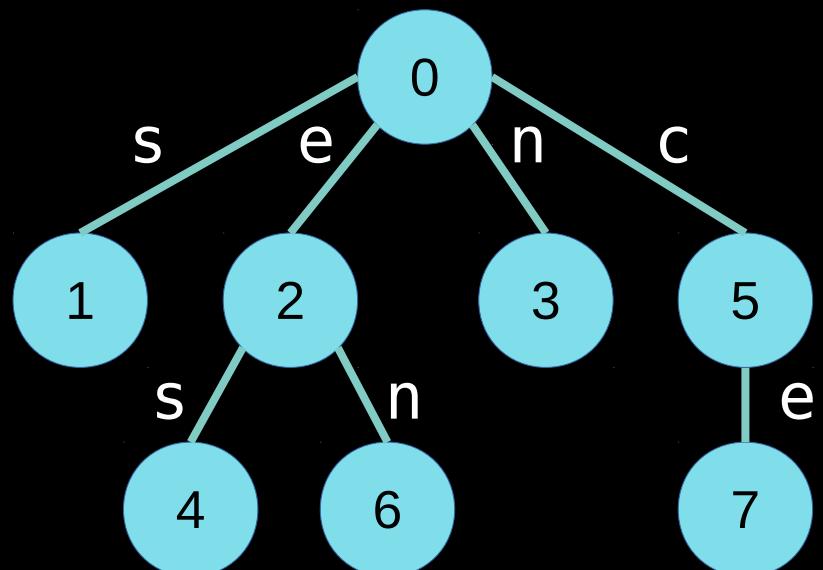


# LZ78 with LZ trie



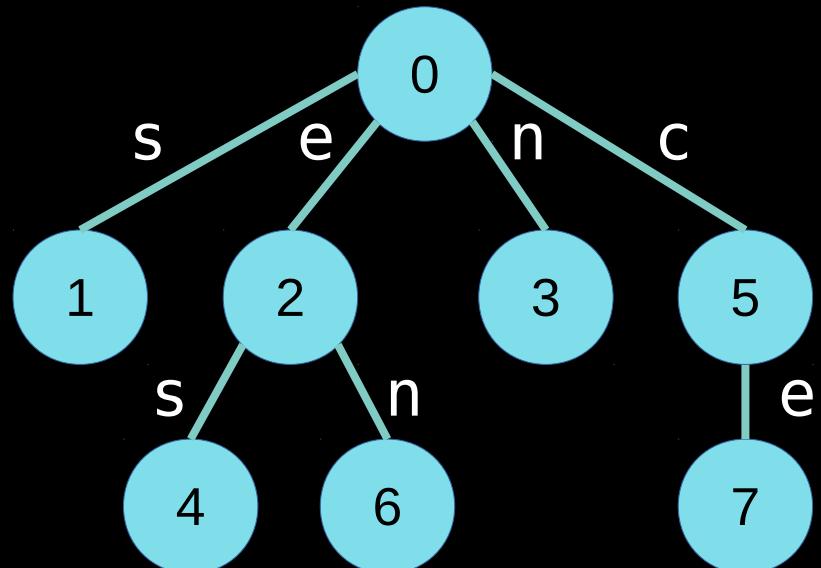
# LZ78 with LZ trie

senescence



# LZ78 with LZ trie

senescence



## Setting

- in what time / memory build LZ trie *practically*?
- $\sigma$ : alphabet size
- $z$ : number of factors

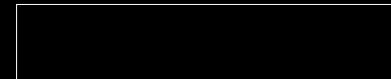
# trie representations

- binary trie [folklore]
- ternary trie [Bentley, Sedgewick'97]
- hash trie
- compact hash trie
- rolling hash trie



new

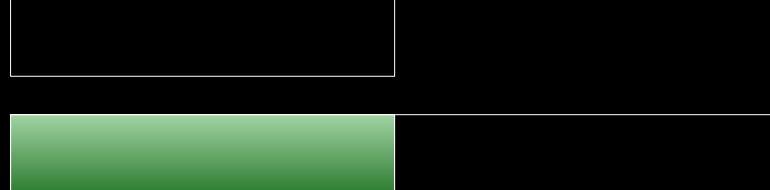
# trie representations

- binary trie [folklore]
  - ternary trie [Bentley, Sedgewick'97]
  - hash trie
  - compact hash trie
  - rolling hash trie
- all use dynamic arrays
- new
- 
- 

# trie representations

- binary trie [folklore]
  - ternary trie [Bentley, Sedgewick'97]
  - hash trie
  - compact hash trie
  - rolling hash trie
- new
- all use dynamic arrays
- double space when full
- 

# trie representations

- binary trie [folklore]
  - ternary trie [Bentley, Sedgewick'97]
  - hash trie
  - compact hash trie
  - rolling hash trie
- new
- all use dynamic arrays
- double space when full
- 

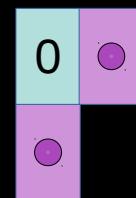
# trie representations

- binary trie [folklore]
  - ternary trie [Bentley, Sedgewick'97]
  - hash trie
  - compact hash trie
  - rolling hash trie
- new      all use dynamic arrays  
double space when full

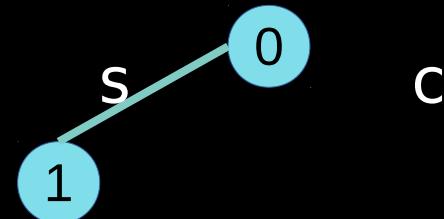


# binary trie

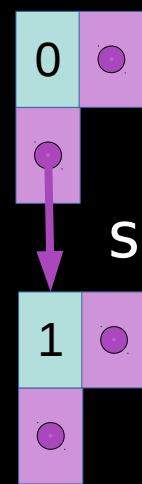
## senescence



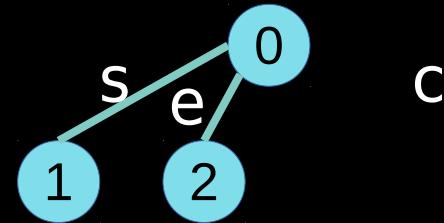
# binary trie



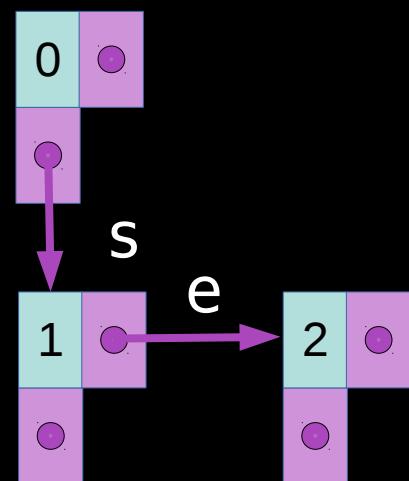
# Senescence



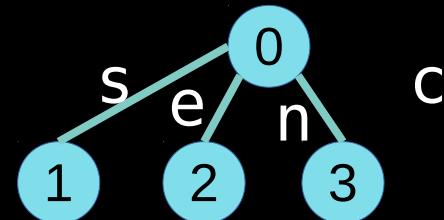
# binary trie



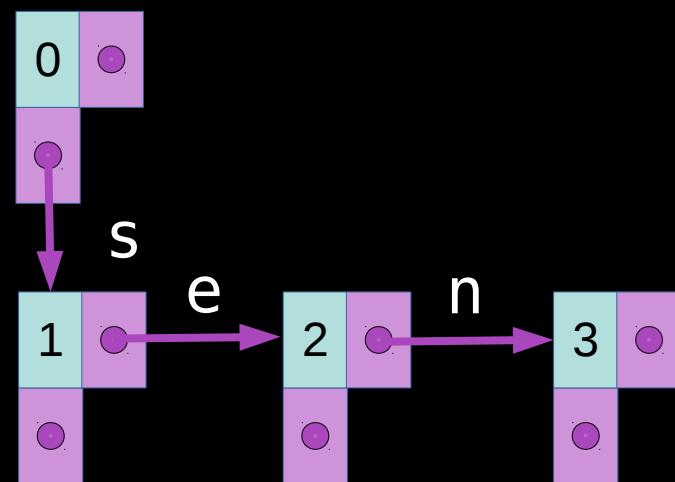
**s**enescence



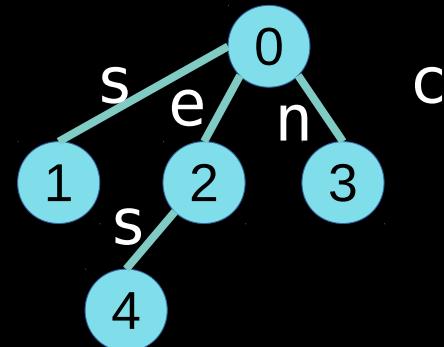
# binary trie



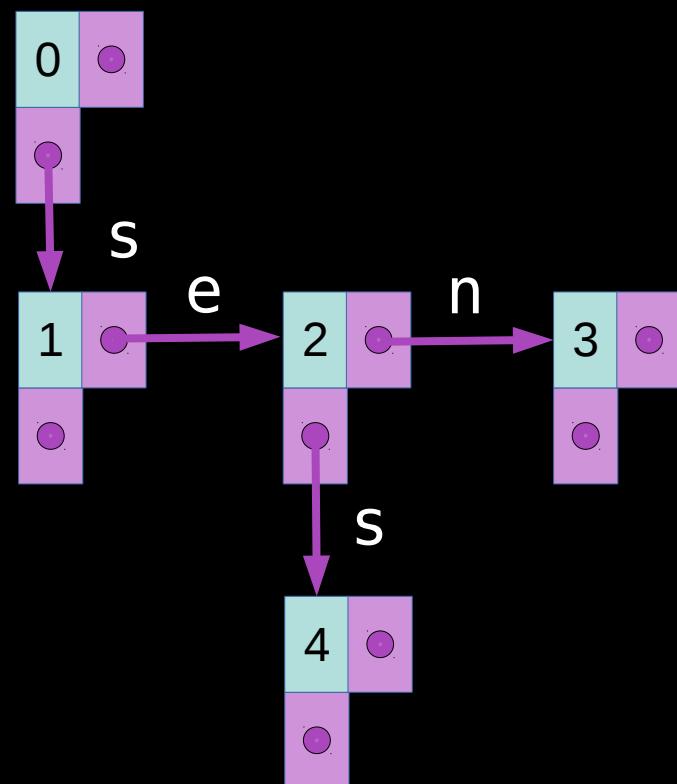
# senescence



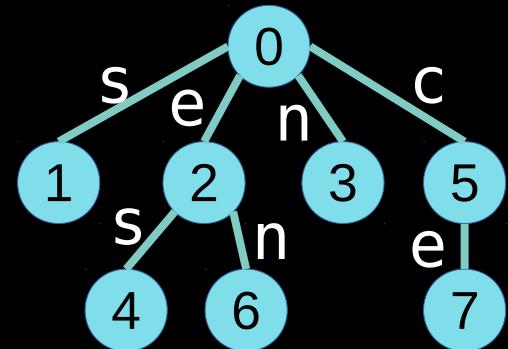
# binary trie



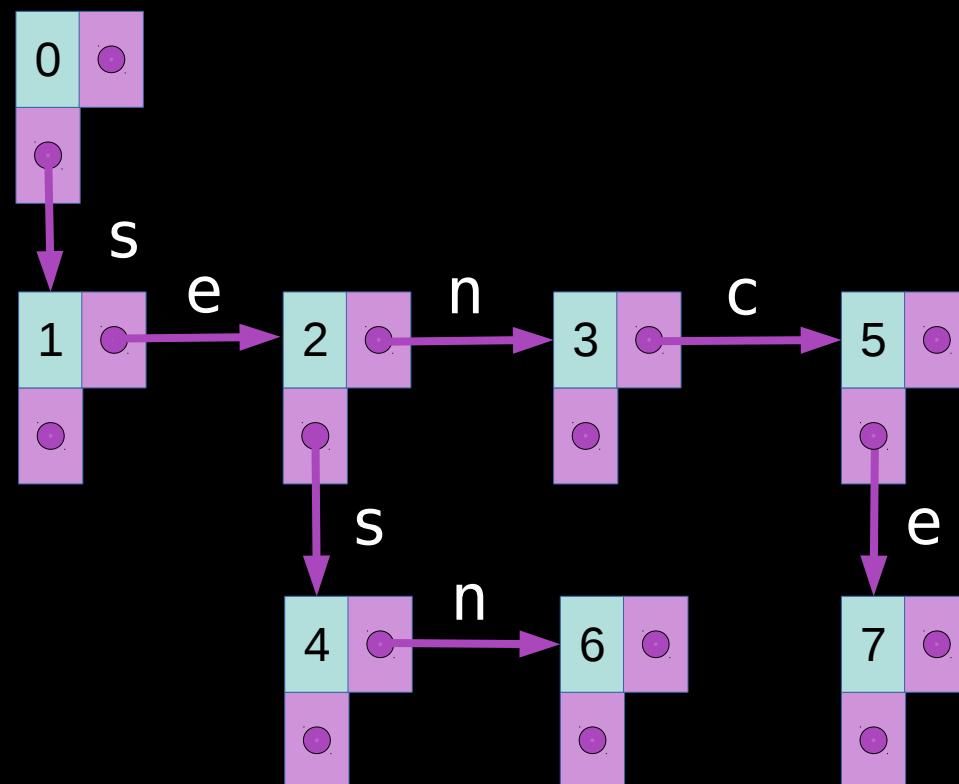
senescence



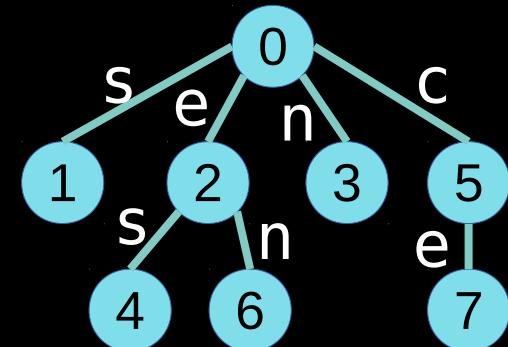
# binary trie



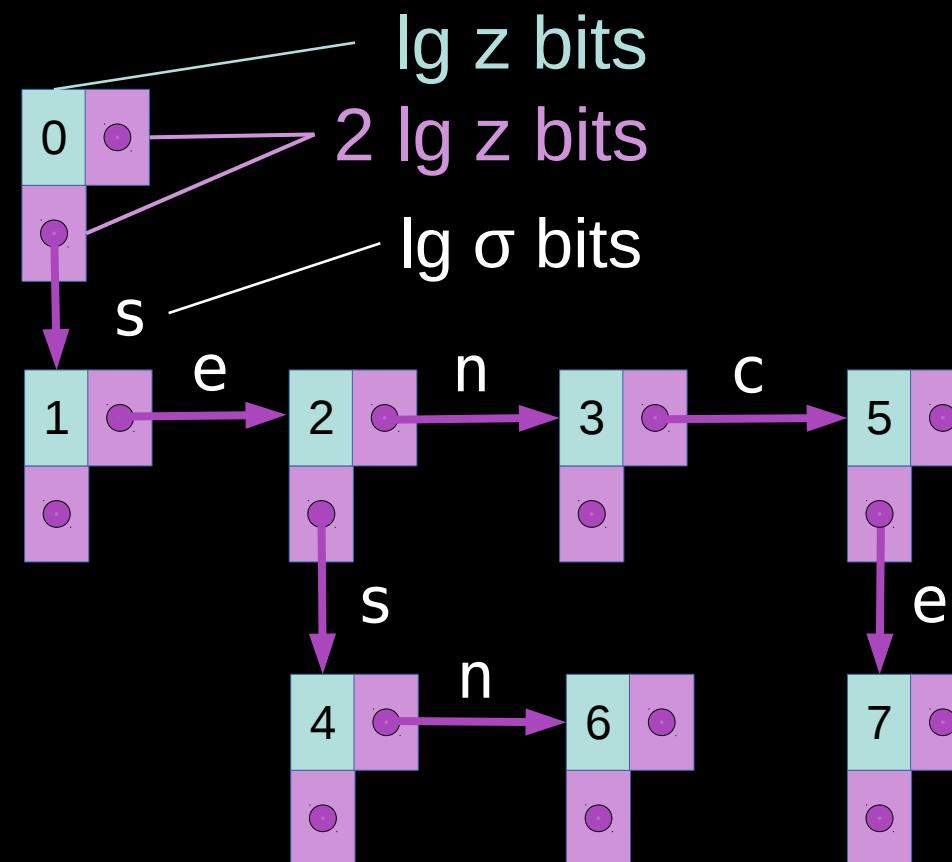
senescence



# binary trie

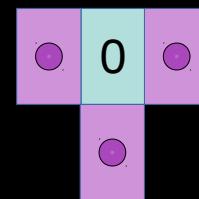


senescence

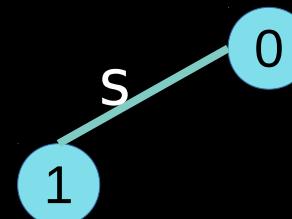


# ternary trie

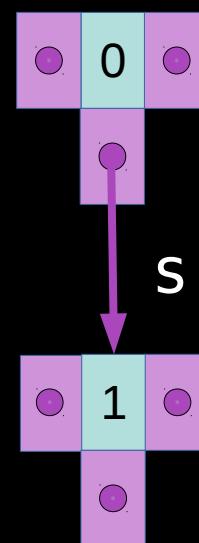
senescence



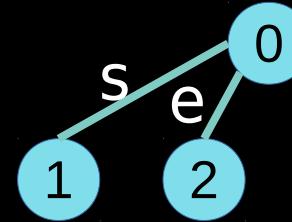
# ternary trie



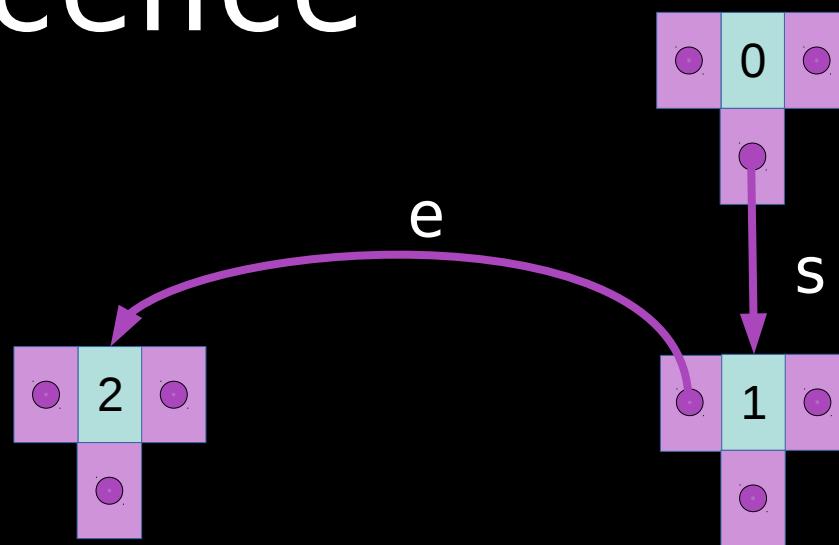
# Senescence



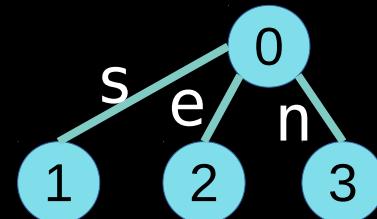
# ternary trie



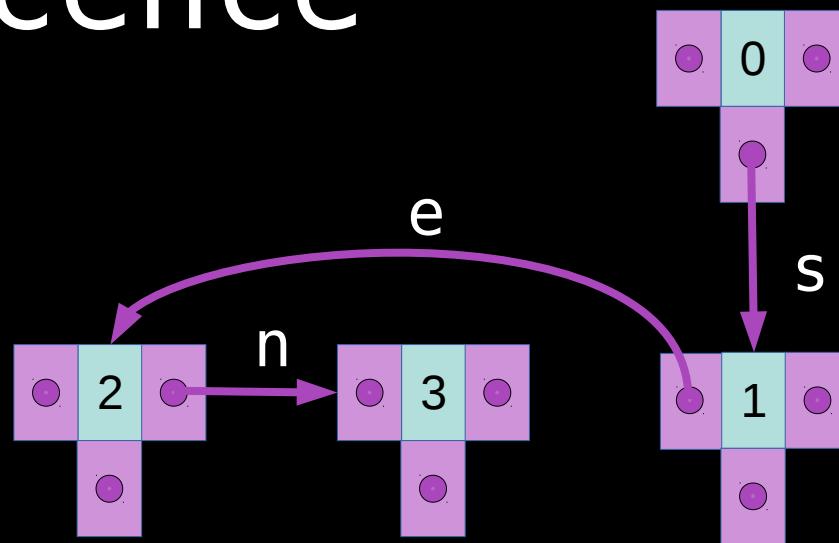
**s**enescence



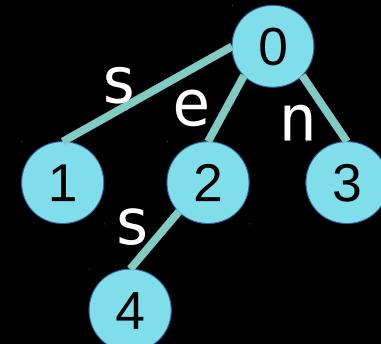
# ternary trie



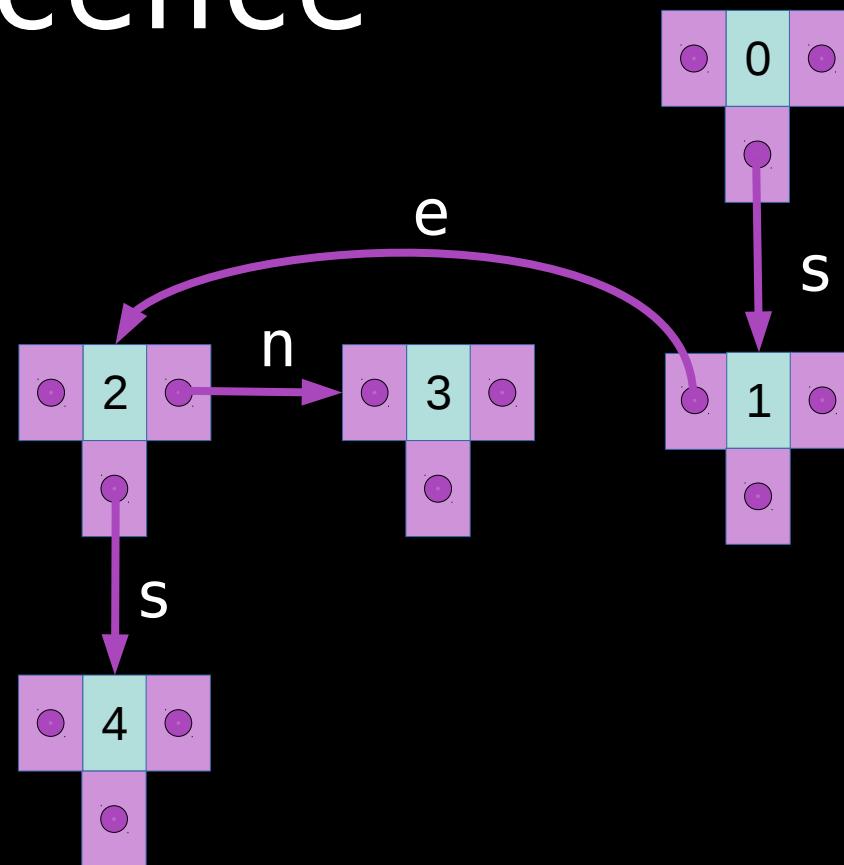
# senescence



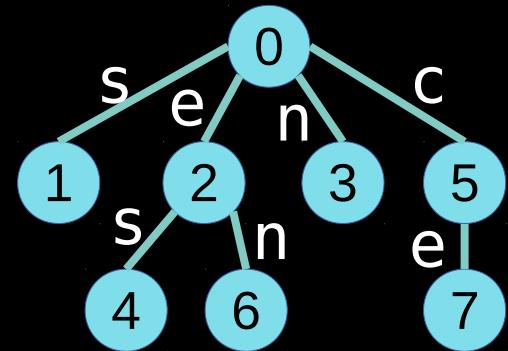
# ternary trie



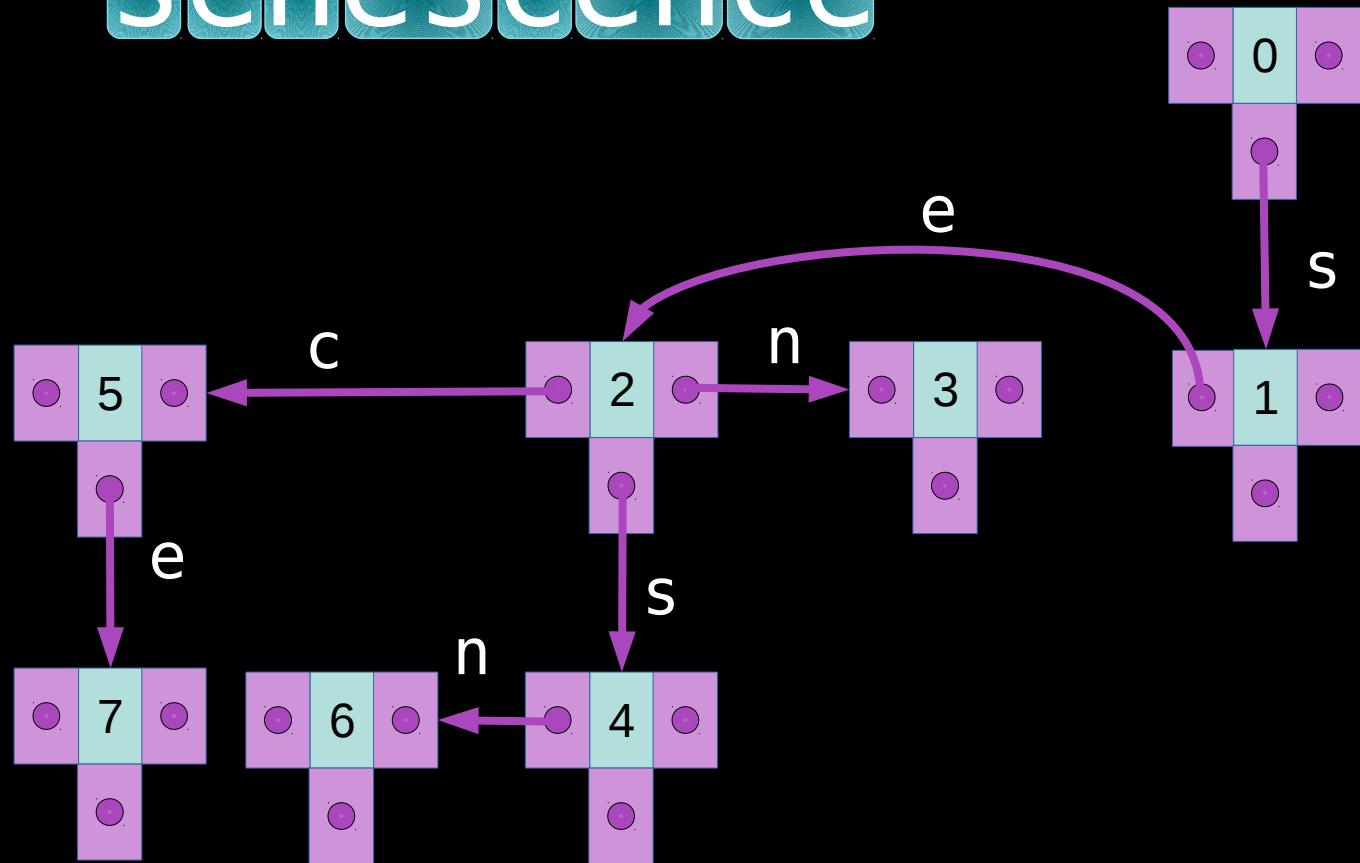
senescence



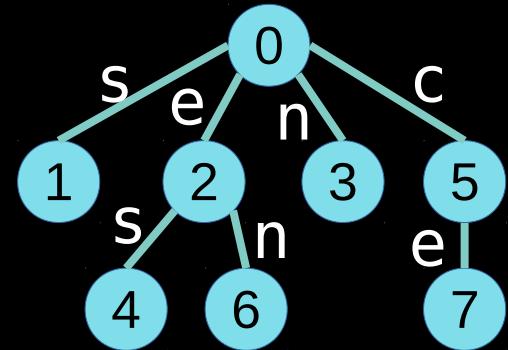
# ternary trie



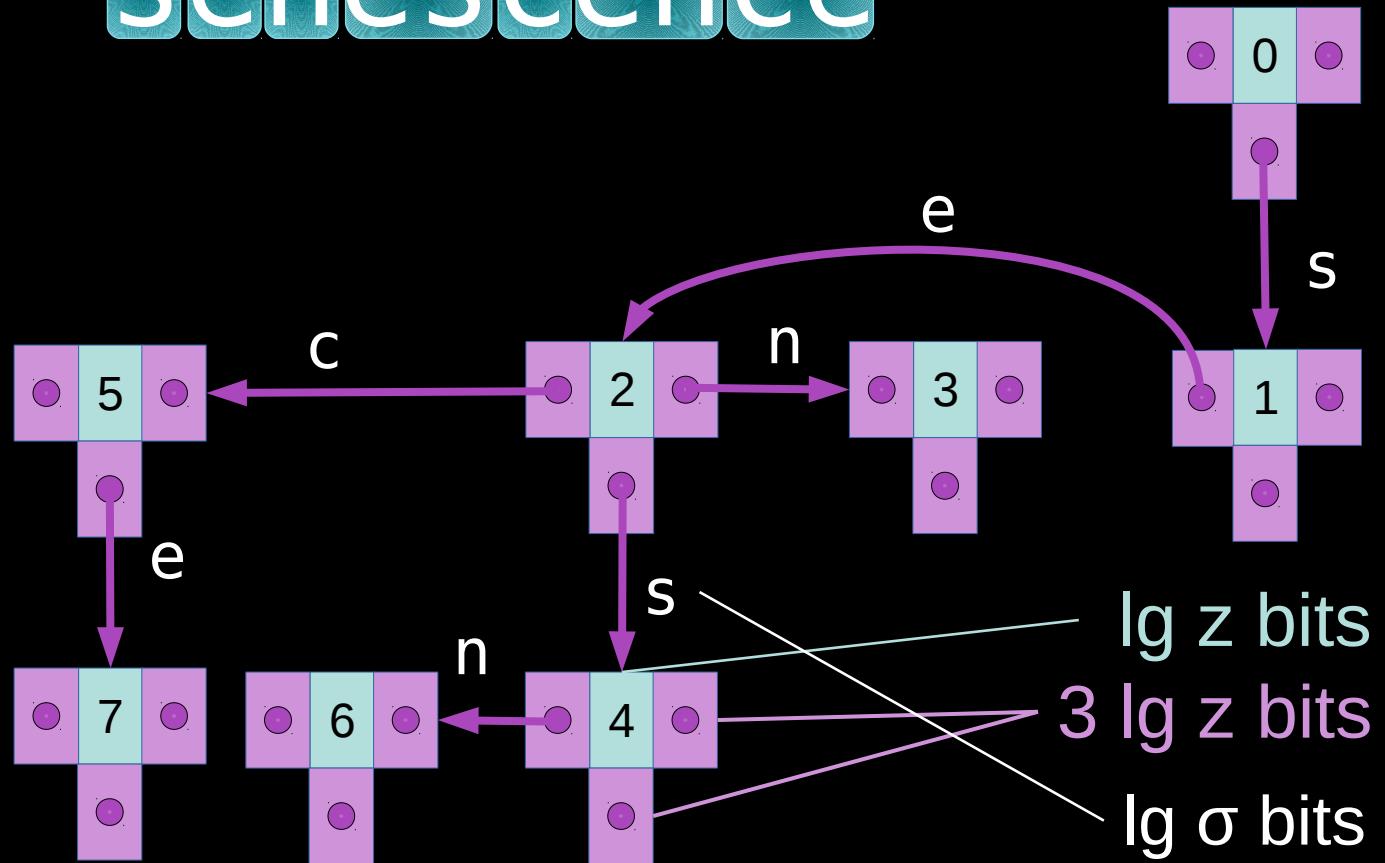
senescence



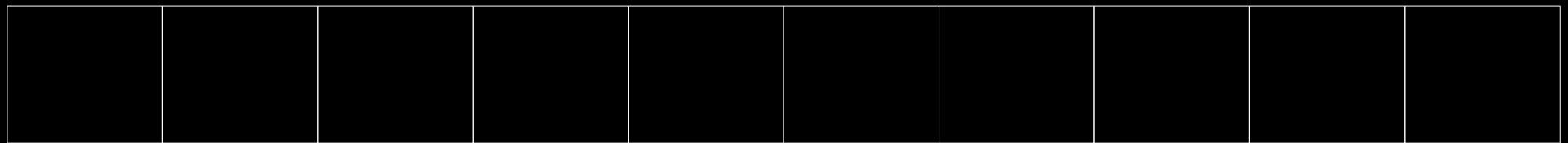
# ternary trie



senescence

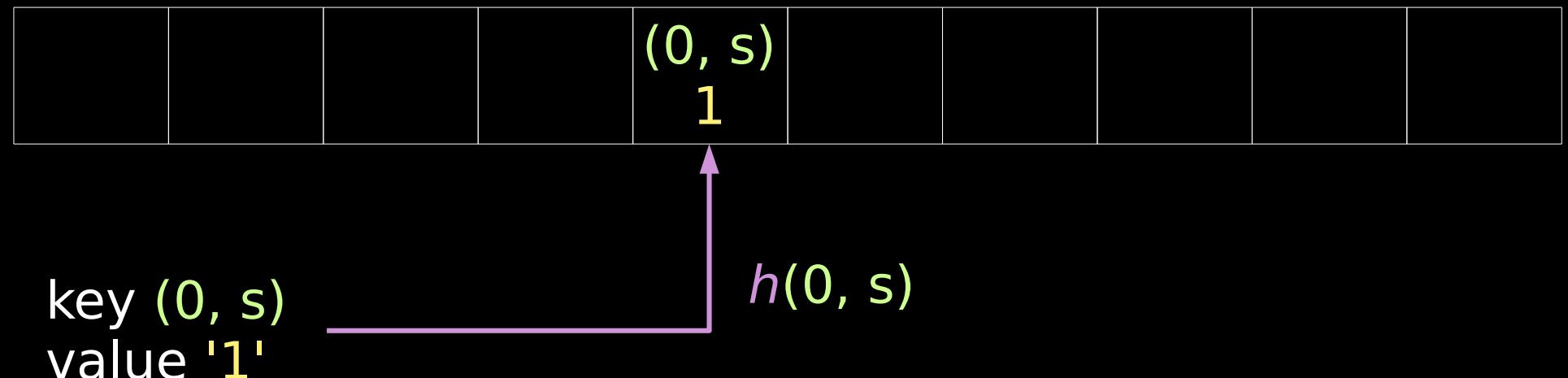
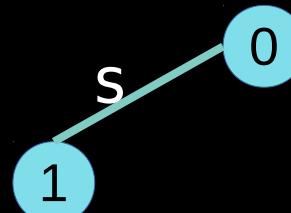


# hash trie



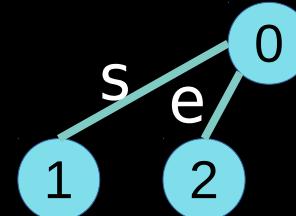
senescence

# hash trie



Senescence

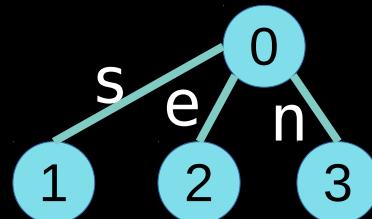
# hash trie



		(0, e) 2	(0, s) 1					
key (0, e) value '2'		$h(0, e)$						

senescence

# hash trie



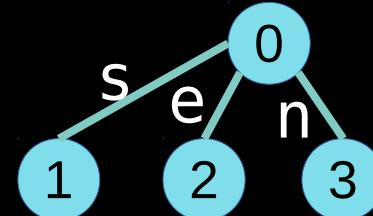
		(0, e) 2	(0, s) 1	(0, n) 3			
--	--	-------------	-------------	-------------	--	--	--

key (0, n)  
value '3'

$h(0, 2)$

senescence

# hash trie



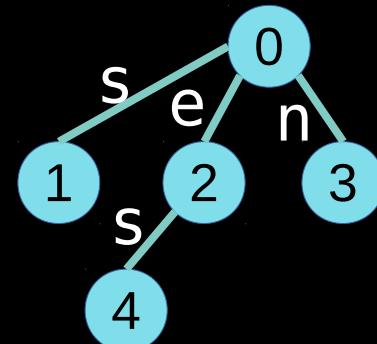
		(0, e) 2	(0, s) 1	(0, n) 3			
--	--	-------------	-------------	-------------	--	--	--

$h(0, e)$

key (0, e)  
value '2'

senescence

# hash trie



		(0, e) 2	(0, s) 1	(0, n) 3	(2, s) 4	
--	--	-------------	-------------	-------------	-------------	--

key (2, s)  
value '4'

senescence

$h(2, s)$

key  $(i, s)$        $\lg z$  bits  
value  $j$        $\lg \sigma$  bits      }  
                     $\lg z$  bits       $2\lg z + \lg \sigma$  bits

		$(0, e)$ 2		$(0, s)$ 1		$(0, n)$ 3		$(2, e)$ 4	
--	--	---------------	--	---------------	--	---------------	--	---------------	--

key  $(i, s)$        $\lg z$  bits  
 value  $j$        $\lg \sigma$  bits      }  
                    $\lg z$  bits      }       $2\lg z + \lg \sigma$  bits

		$(0, e)$ 2		$(0, s)$ 1		$(0, n)$ 3		$(2, e)$ 4	
--	--	---------------	--	---------------	--	---------------	--	---------------	--

M

- table size  $M$
- load factor  $\alpha$
- $z \leq \alpha M$

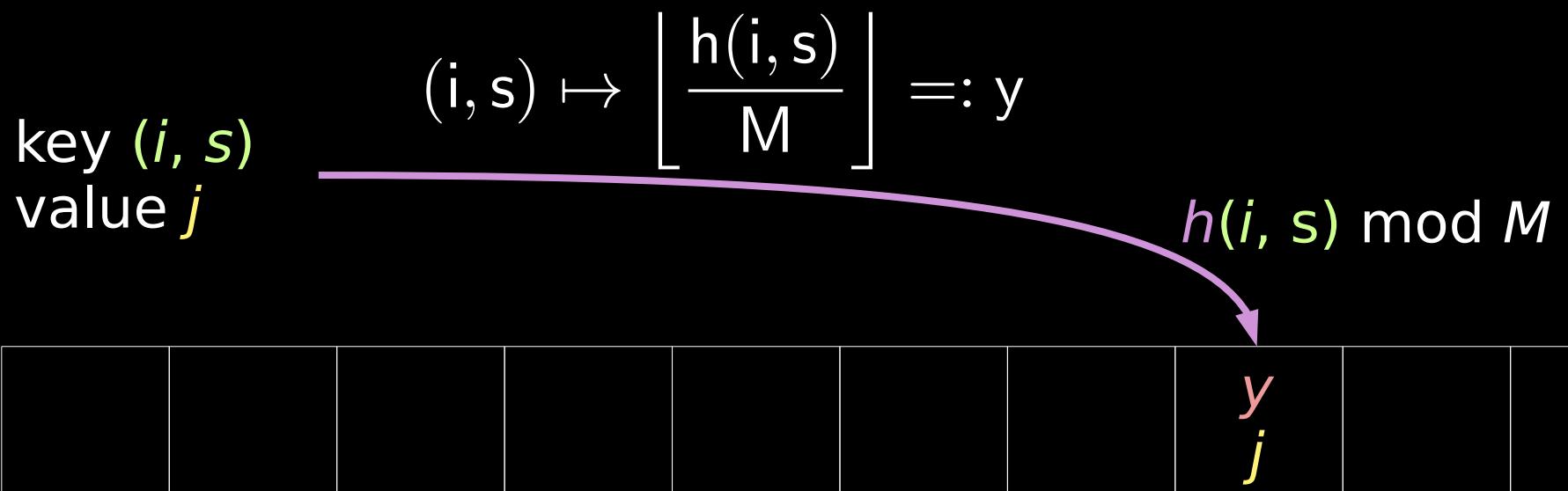
$$\geq \frac{z}{\alpha} (2 \lg z + \lg \sigma) \text{ bits}$$

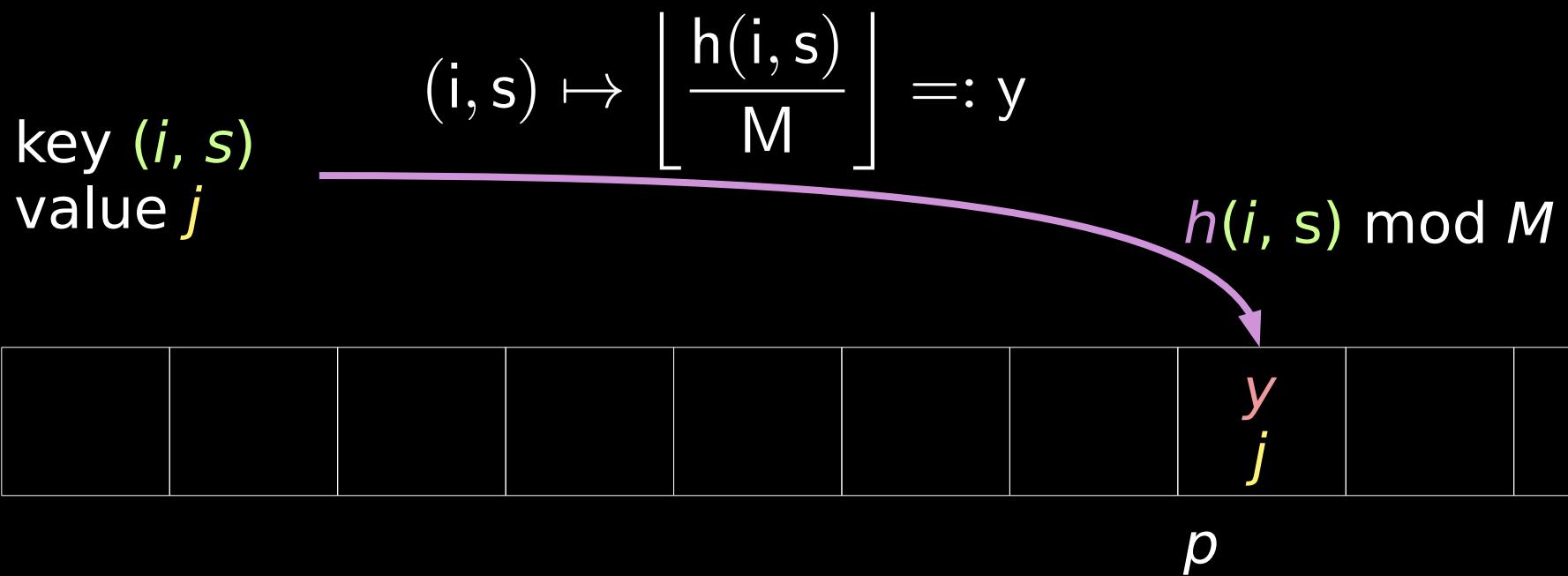
# compact hash trie

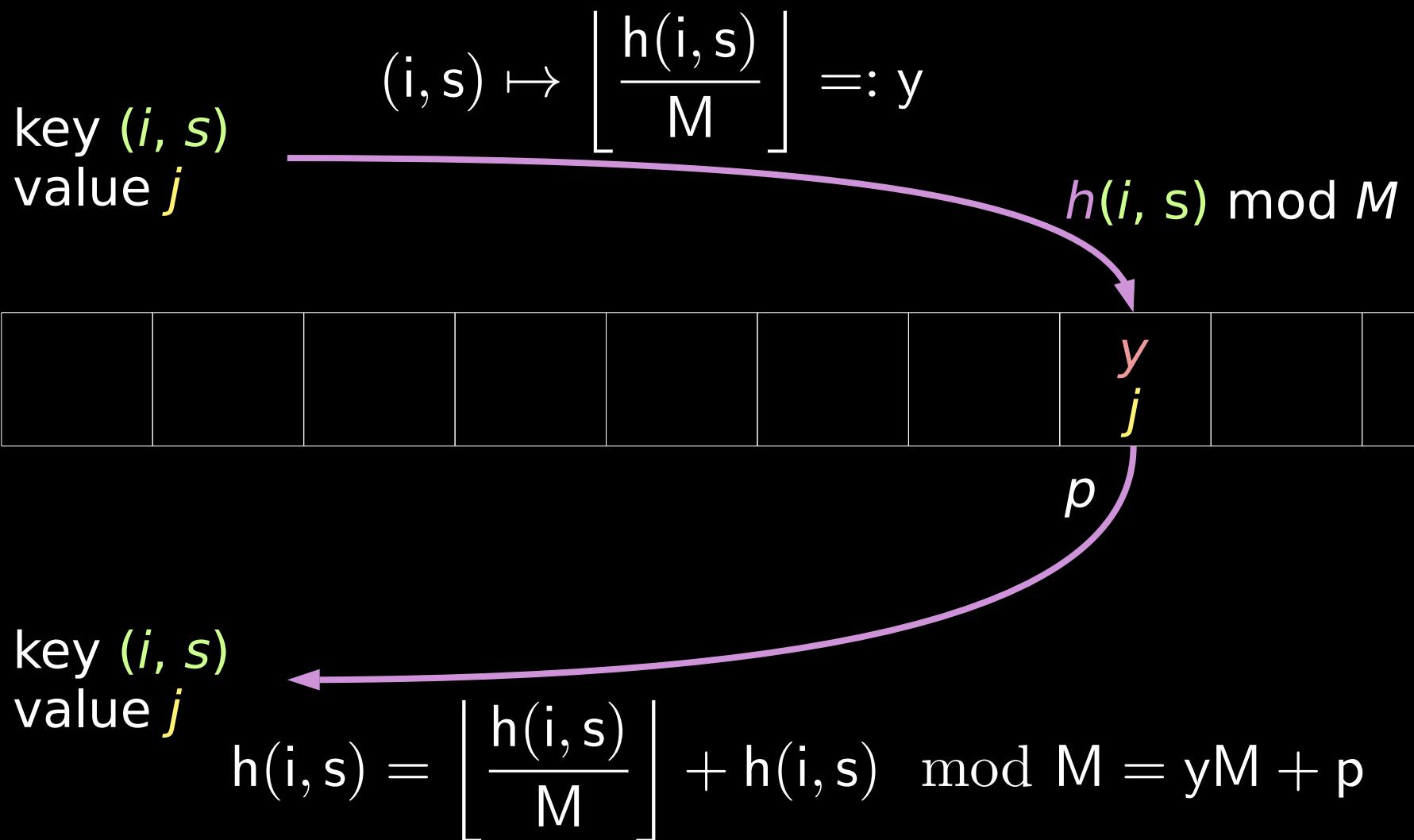
- bijective hash function
  - $(ax + b) \bmod q$
  - xorshift  $x \rightarrow x \oplus (x \gg j)$
- choose  $j, q$  so  $h : [0..\sigma M] \rightarrow [0..\sigma M]$  is bijective:
  - invariant:  $z \leq \alpha M$
  - $|[0..z-1] \times [0..\sigma-1]| = z\sigma \leq \alpha\sigma M \leq \sigma M$

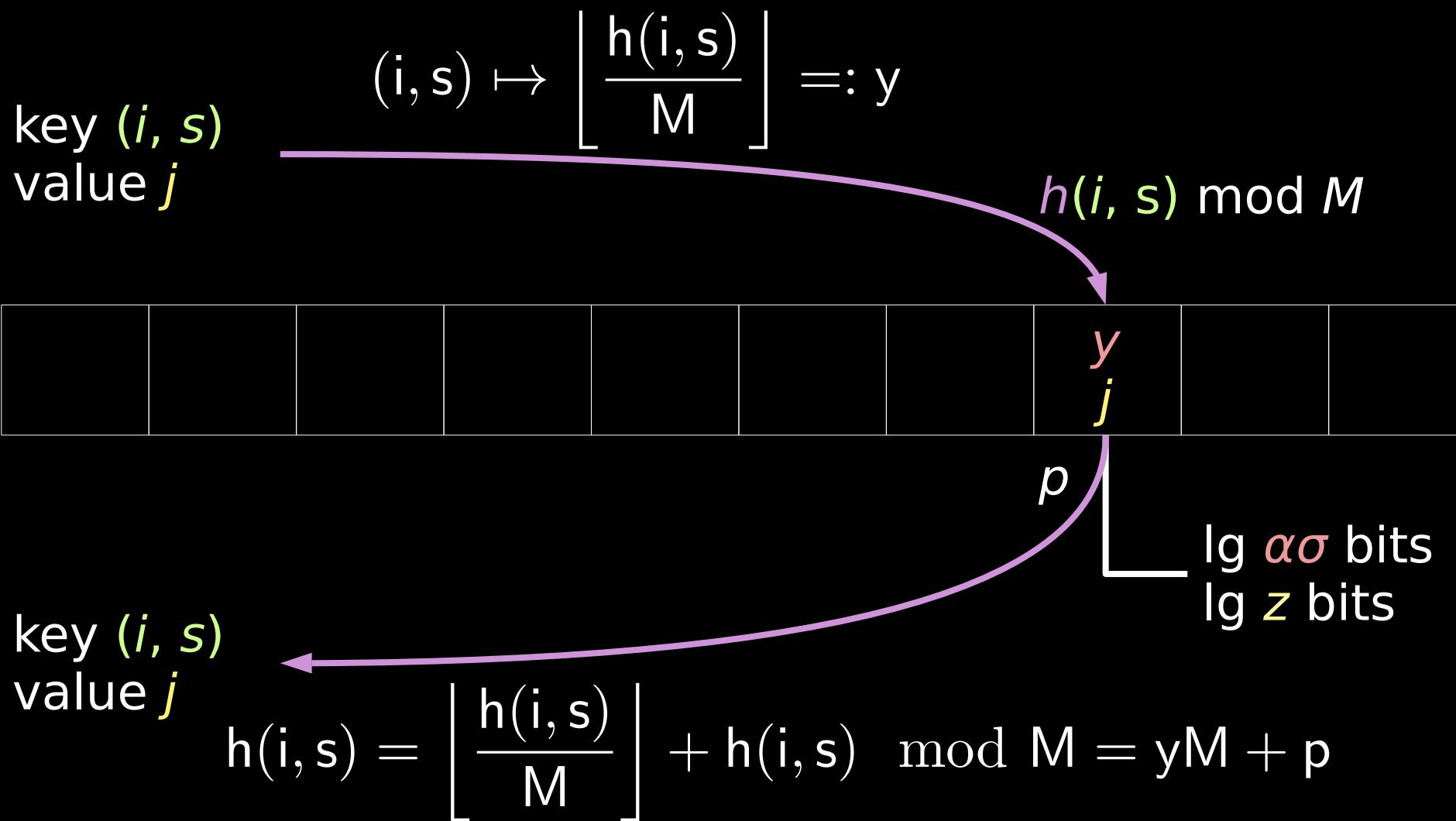
based on [Poyias and Raman'17]

key  $(i, s)$   
value  $j$



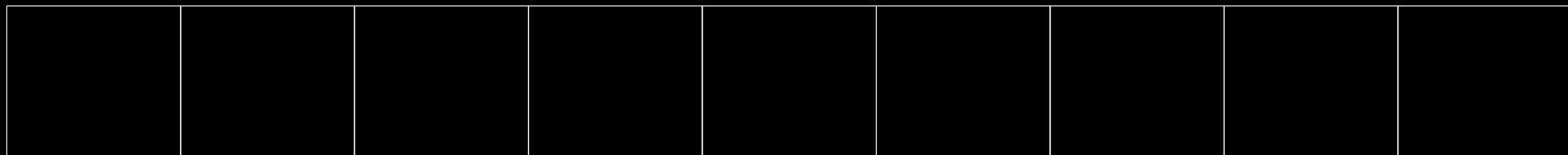




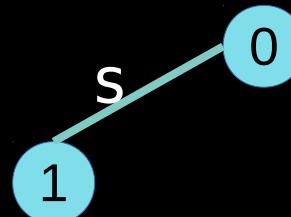


# rolling hash trie

## senescence



# rolling hash trie

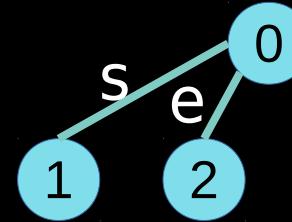


senescence

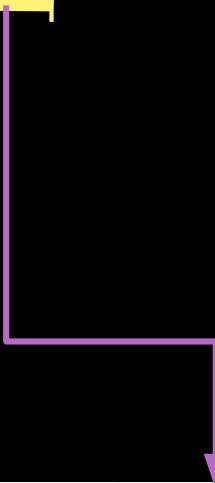
```
[senescence]
```

	$h(s)$ 1							
--	-------------	--	--	--	--	--	--	--

# rolling hash trie

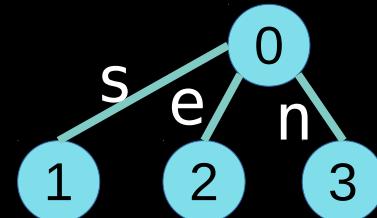


senescence

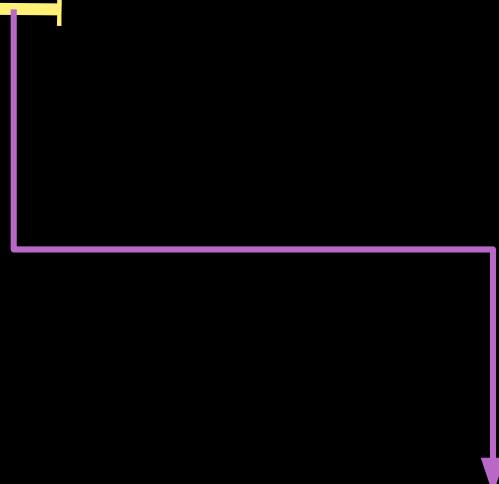


	$h(s)$ 1		$h(e)$ 2					
--	-------------	--	-------------	--	--	--	--	--

# rolling hash trie

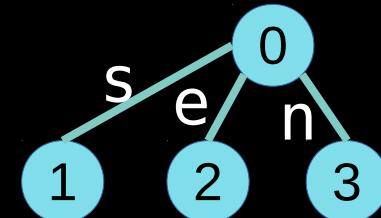


senescence



	$h(s)$ 1		$h(e)$ 2		$h(n)$ 3			
--	-------------	--	-------------	--	-------------	--	--	--

# rolling hash trie

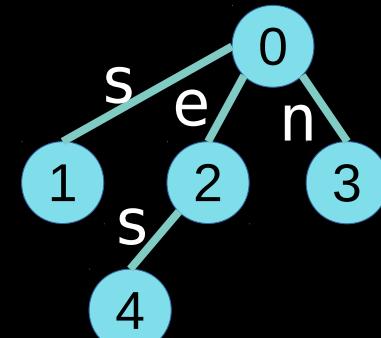


senescence



	$h(s)$ 1		$h(e)$ 2		$h(n)$ 3			
--	-------------	--	-------------	--	-------------	--	--	--

# rolling hash trie

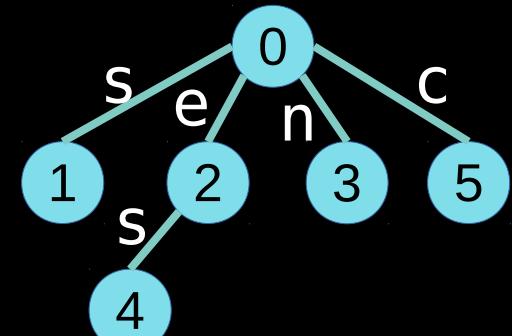


senescence

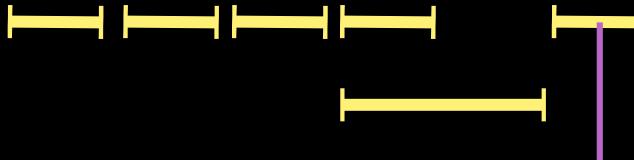


	$h(s)$ 1		$h(e)$ 2		$h(n)$ 3		$h(es)$ 4	
--	-------------	--	-------------	--	-------------	--	--------------	--

# rolling hash trie



senescence



$$h(c) = h(es)$$



	$h(s)$ 1		$h(e)$ 2		$h(n)$ 3		$h(es)$ 4	
--	-------------	--	-------------	--	-------------	--	--------------	--

# properties

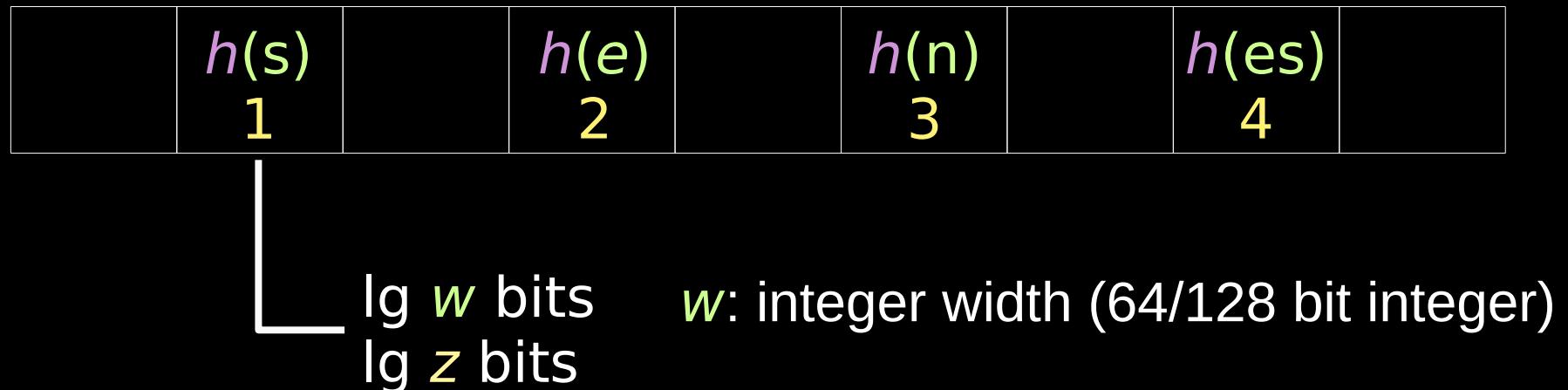
- Monte Carlo algorithm
- store fingerprint instead of keys

	$h(s)$ 1		$h(e)$ 2		$h(n)$ 3		$h(es)$ 4	
--	-------------	--	-------------	--	-------------	--	--------------	--

$w$ : integer width (64/128 bit integer)

# properties

- Monte Carlo algorithm
- store fingerprint instead of keys

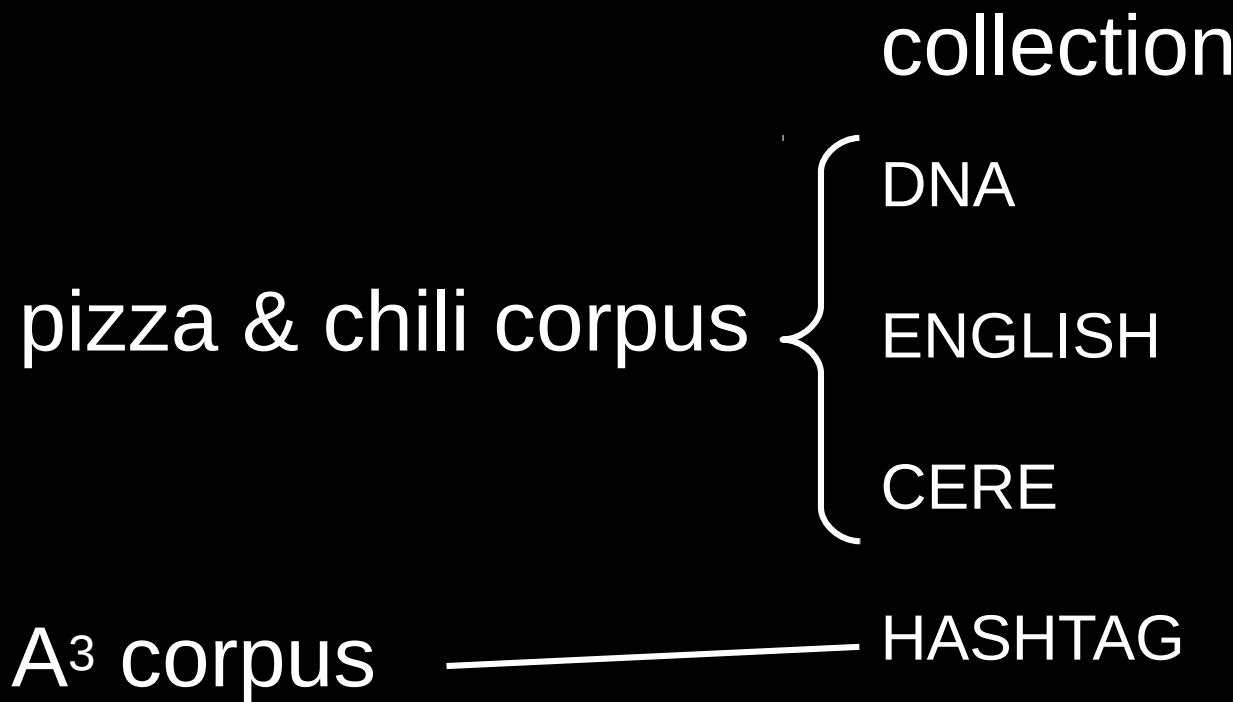


# space

data structure	bits
binary trie	$z(2 \lg z + \lg \sigma)$
ternary trie	$z(3 \lg z + \lg \sigma)$
hash trie	$z/\alpha(2 \lg z + \lg \sigma)$
compact hash trie	$z/\alpha \lg \alpha z \sigma$
rolling hash trie	$z/\alpha \lg w z$

# datasets

200 MiB text files



# technical details

## **type sizes**

- byte alphabet  
 $(\sigma = 256)$
- $w = 64$  bits

# technical details

## type sizes

- byte alphabet  
 $(\sigma = 256)$
- $w = 64$  bits

## hashing

- linear probing
- hash functions:
  - splitmix64 (Vigna'15)
  - rolling:

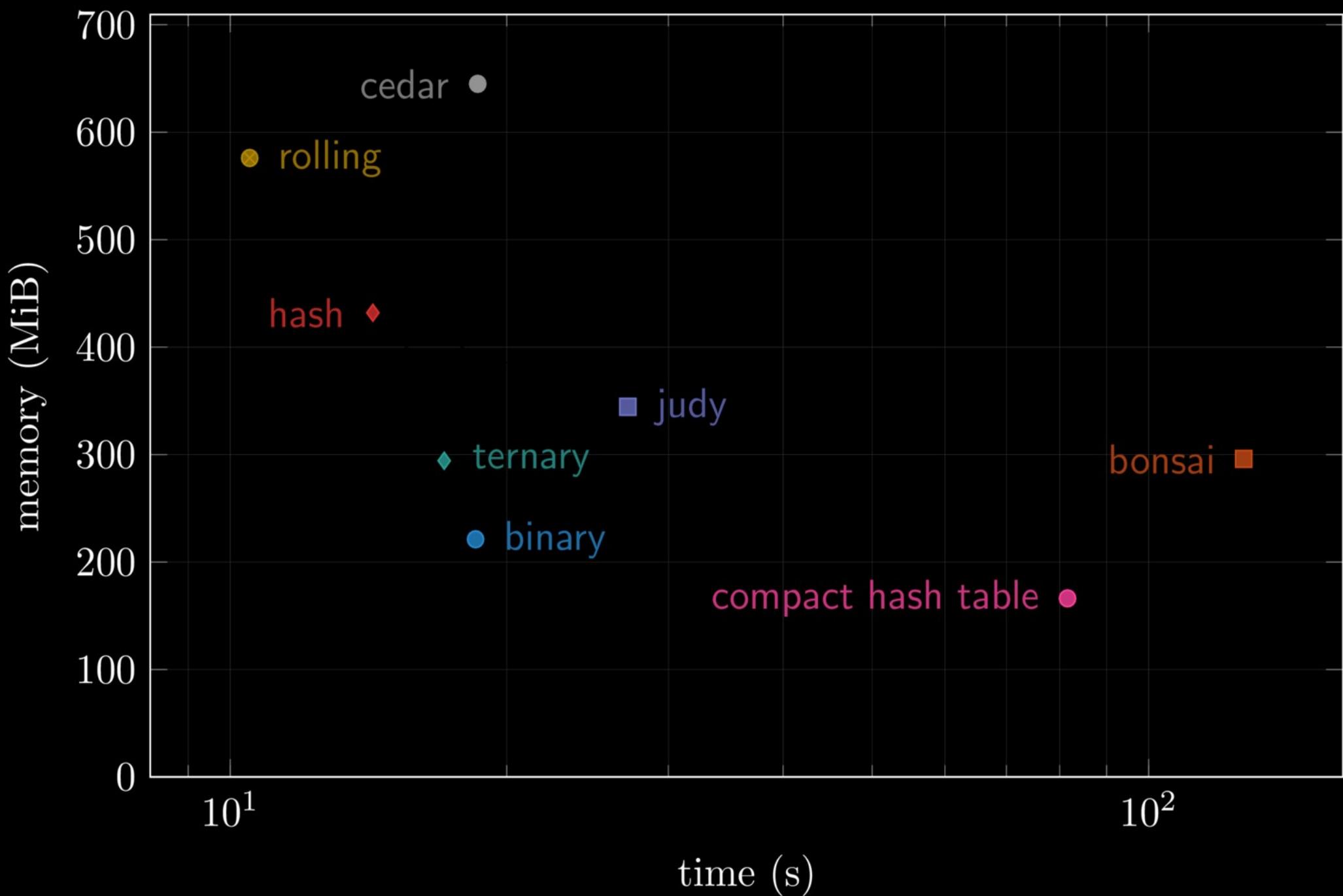
$$h(T) := \sum_{i=1}^{|T|} T[i](\sigma + 1)^{|T|-i} \mod 2^{64}$$

$\sigma+1 = 257$  prime number

# competitors

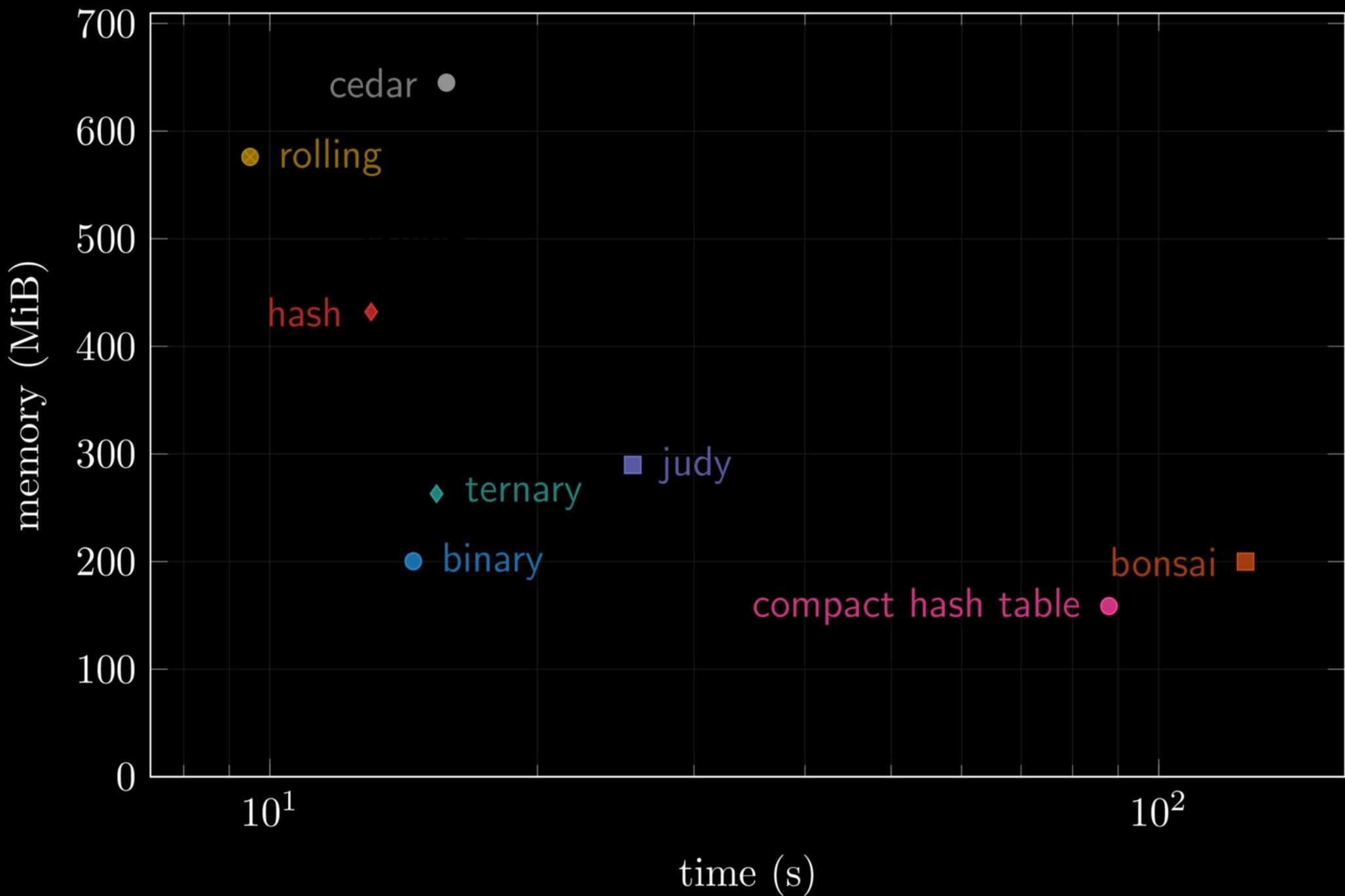
- judy array [Baskins'02]
  - 256-ary radix tries
- cedar [Yoshinaga and Kitsuregawa'14]
  - double array variant
- bonsai [Poyias and Raman'15]
  - trie with compact hashing

# HASHTAG



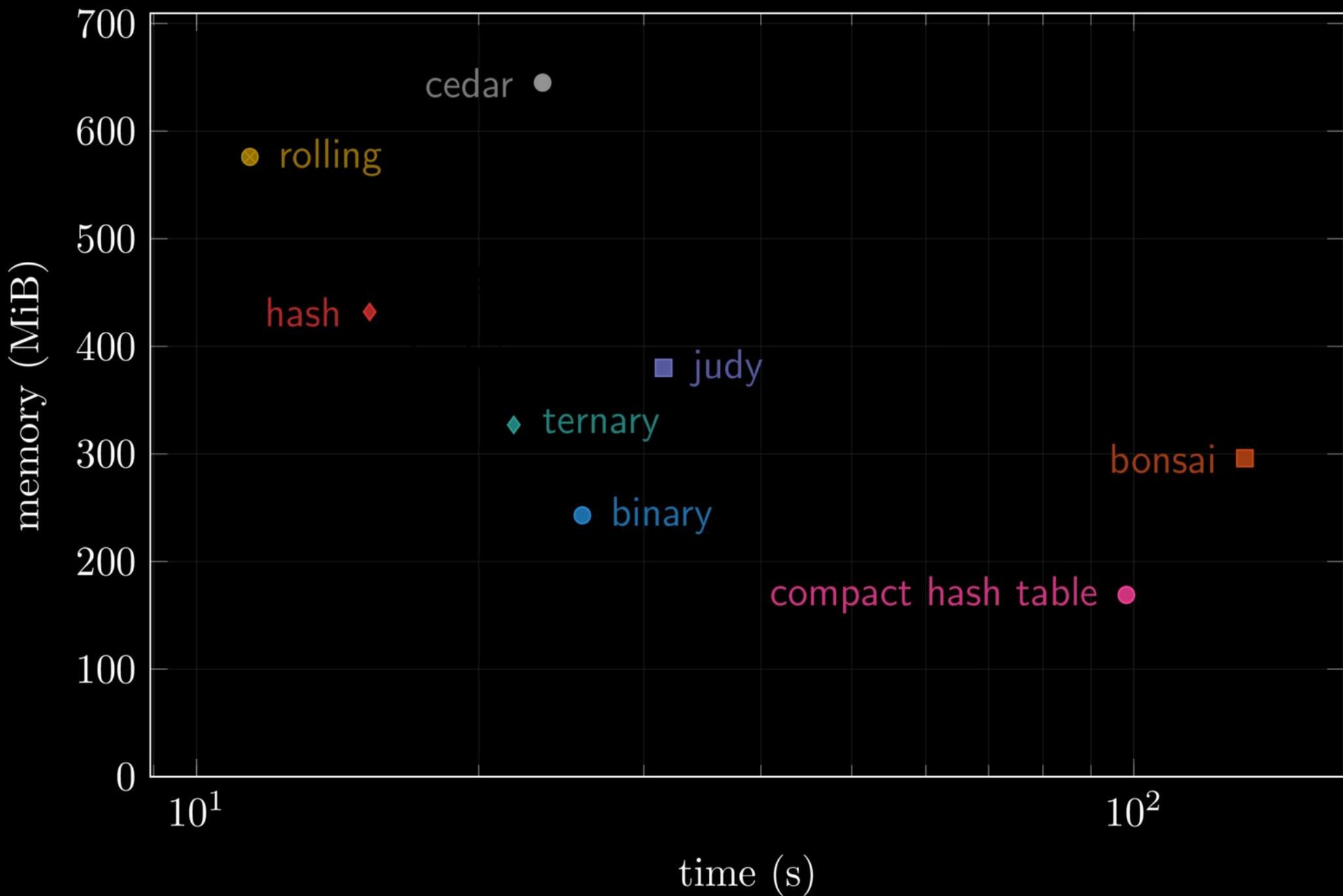
# PC-DNA

58



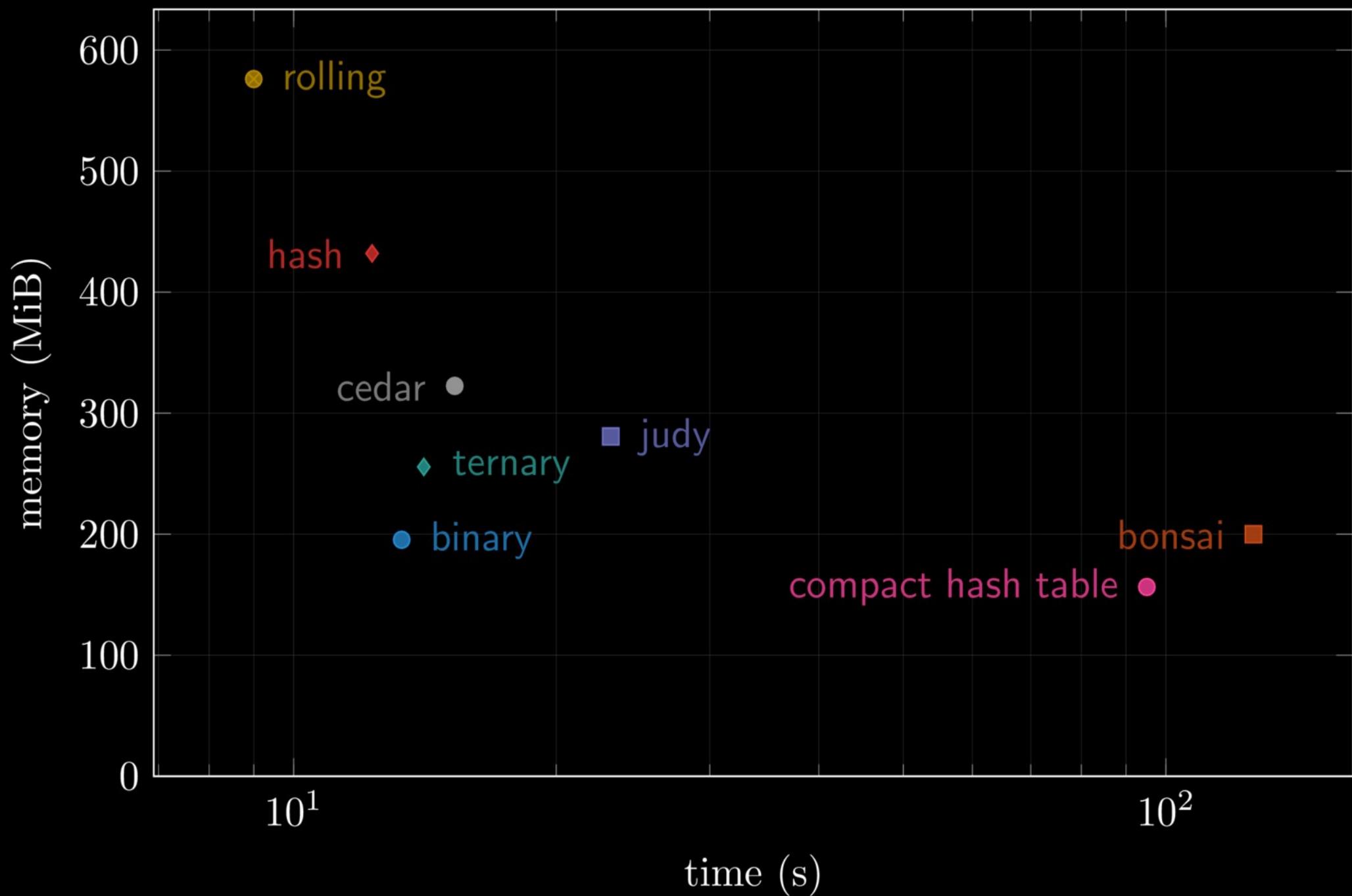
# PC-ENGLISH

59



# PCR-CERE

60



# conclusion

- **rolling** fastest, but output maybe corrupt
- **hash** fast + correct
- **ternary** best trade-off speed  $\leftrightarrow$  memory
- **binary** slower, but less memory
- **compact hash trie** memory friendliest
- competitors not in Pareto front

# conclusion

- **rolling** fastest, but output maybe corrupt
- **hash** fast + correct
- **ternary** best trade-off speed  $\leftrightarrow$  memory
- **binary** slower, but less memory
- **compact hash trie** memory friendliest
- competitors not in Pareto front

Thank you for your attention. Any questions are welcome!