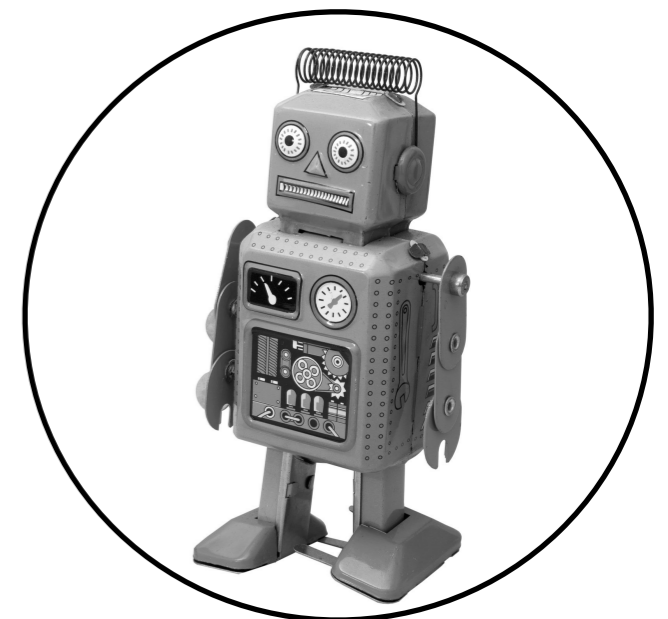


Database Kernels for Data Exploration

Stratos Idreos



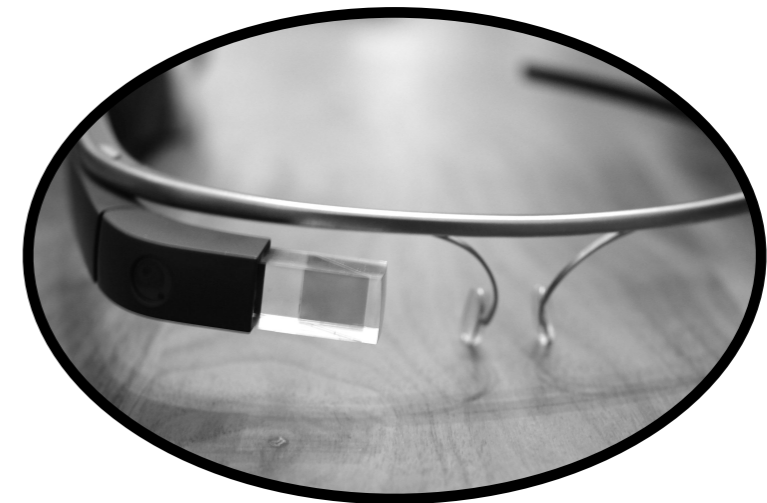
today



today



tomorrow



soon everyone will need to be a “data scientist”

My data is too big :(





data exploration



not always sure
what we are looking for
(until we find it)

data has always been **big**

volume

velocity

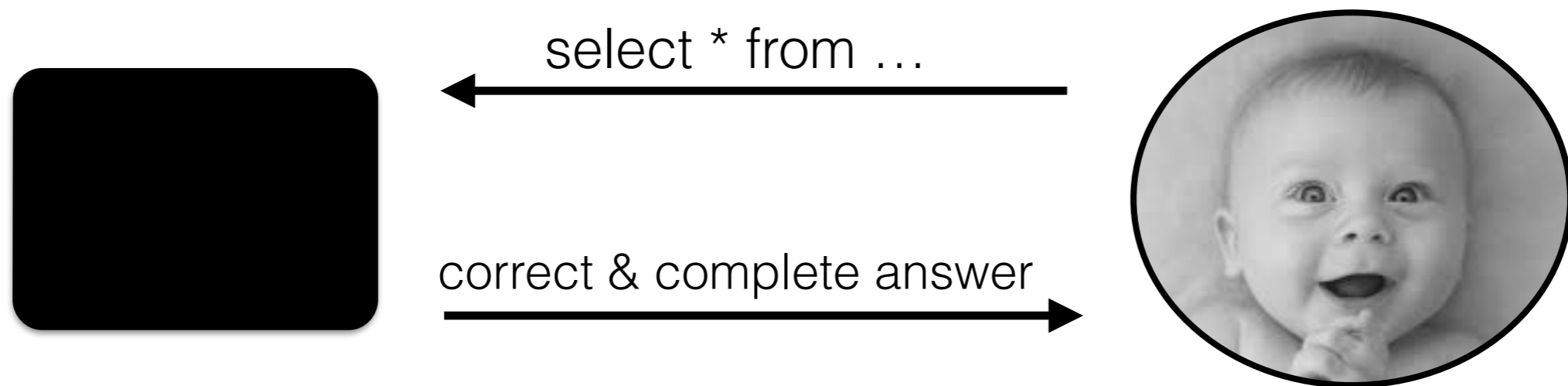
variety

veracity



how future data systems should look like?

the premise of the black box



the perfect data system

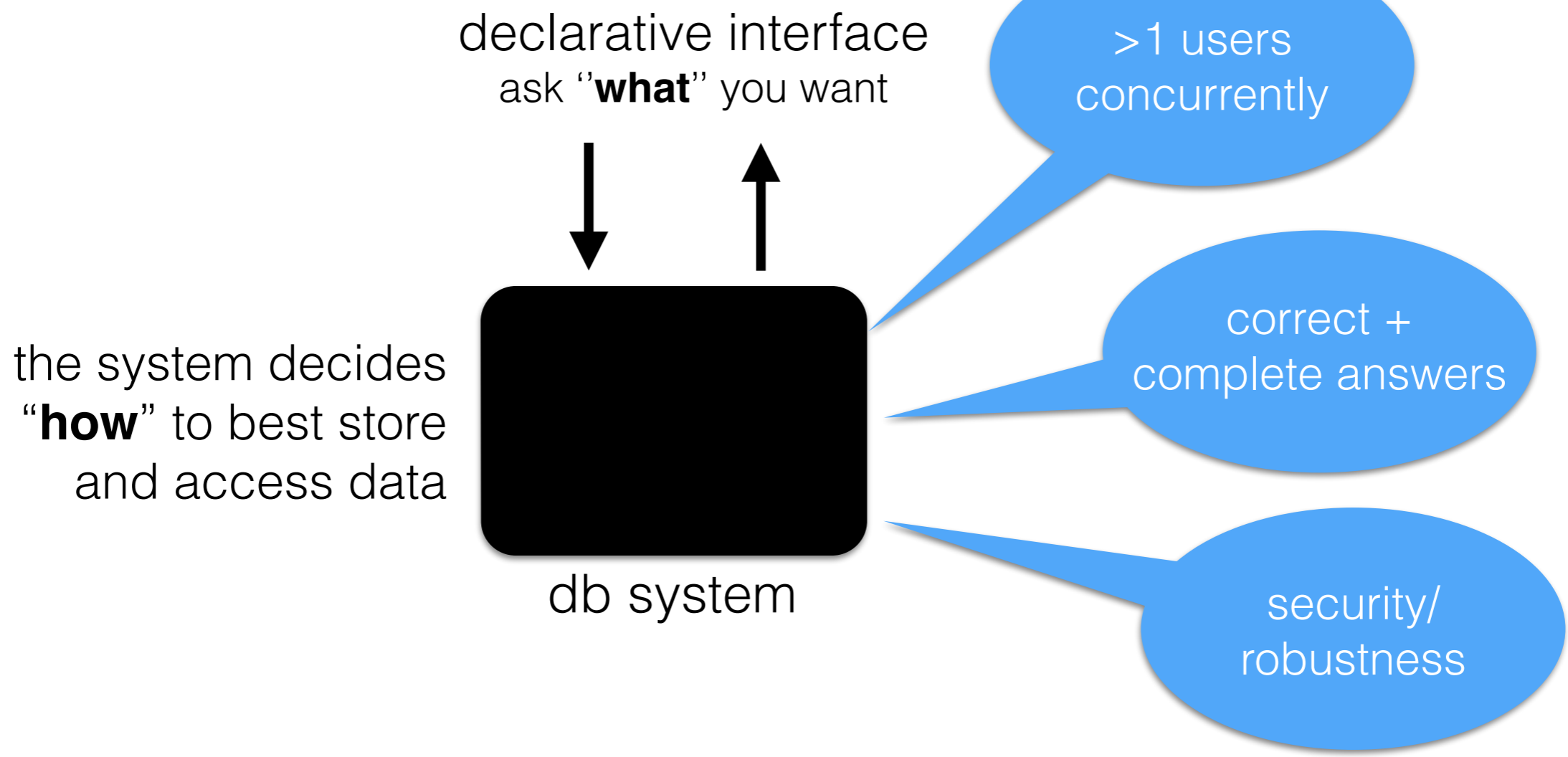
declarative interface
ask "**what**" you want

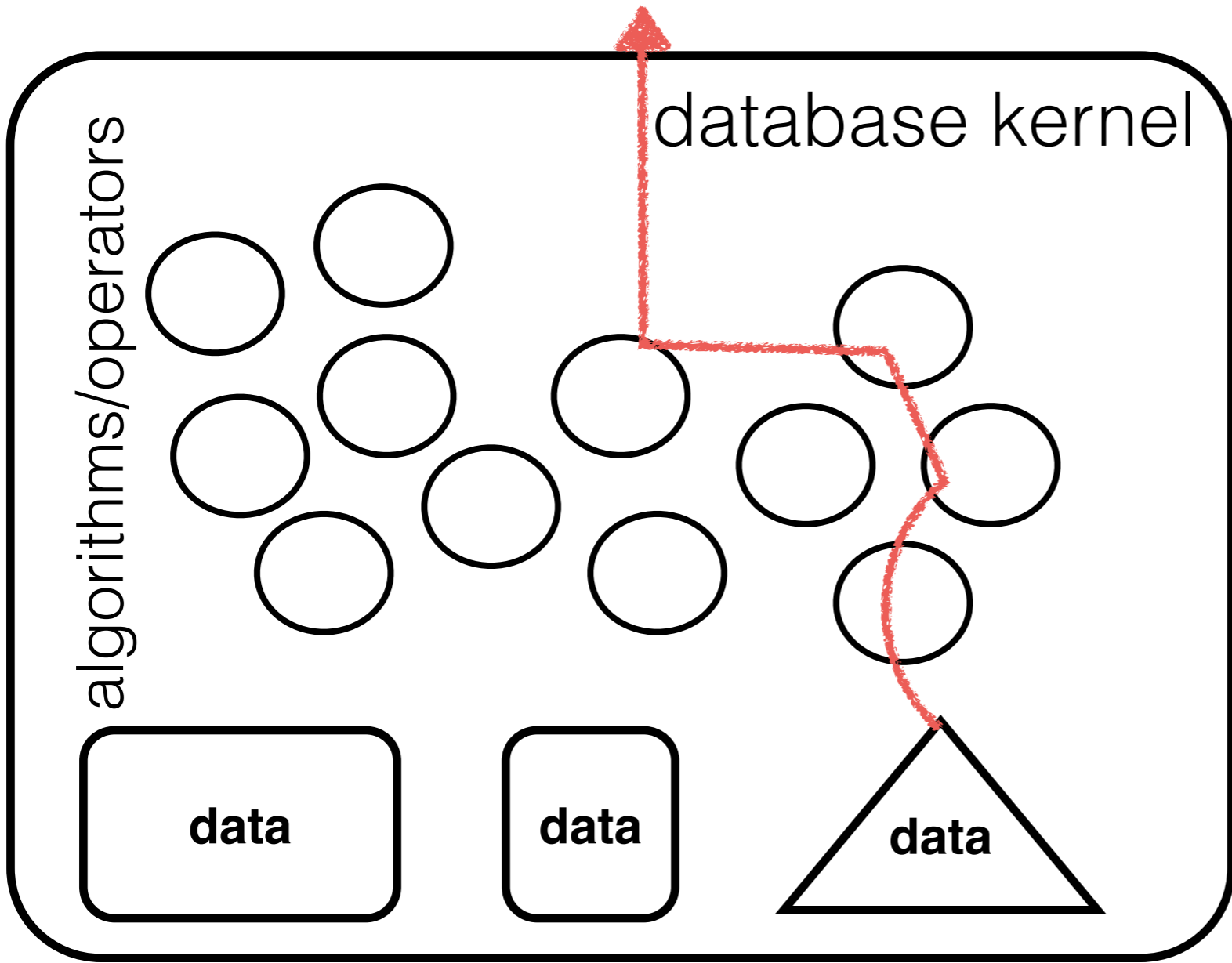


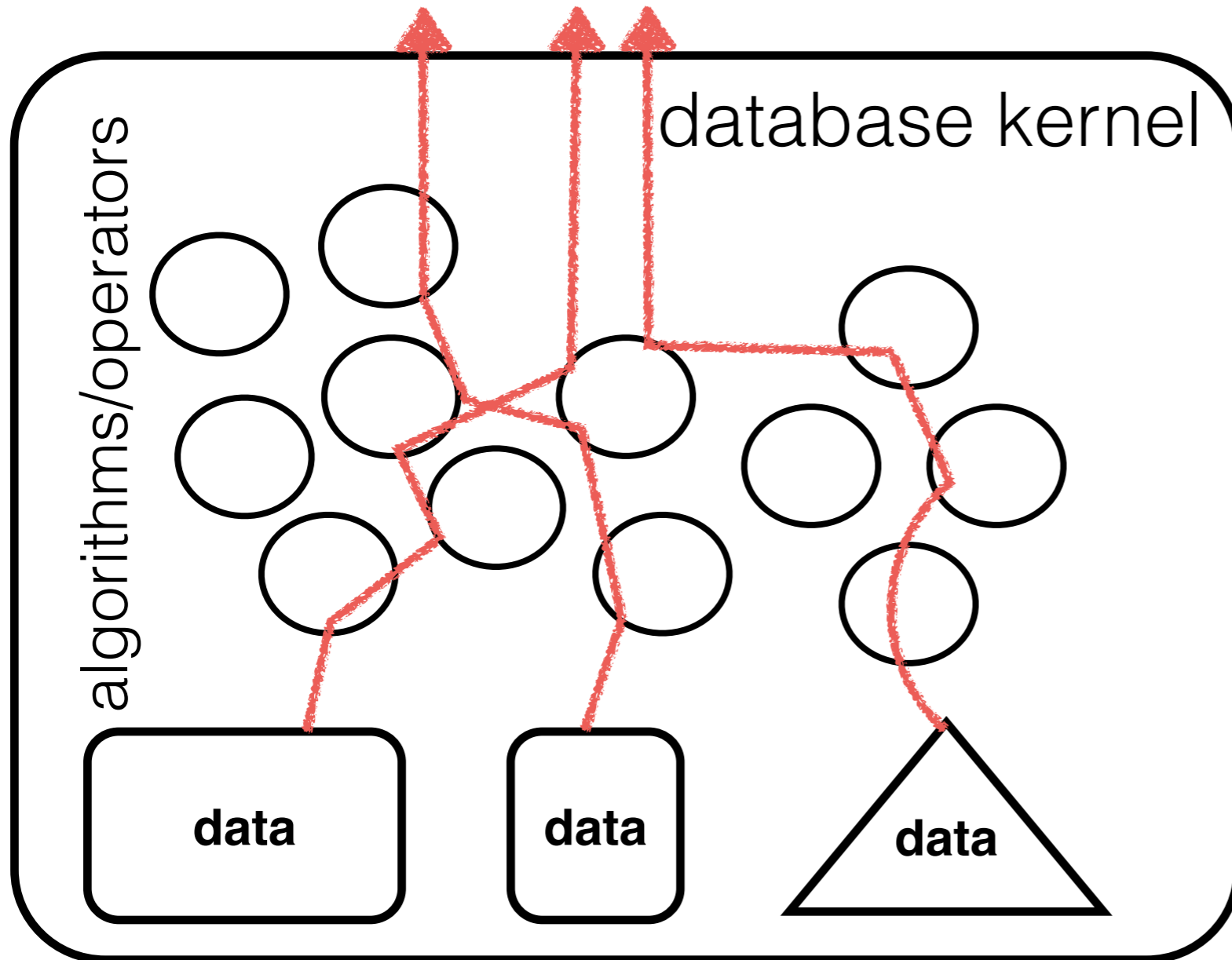
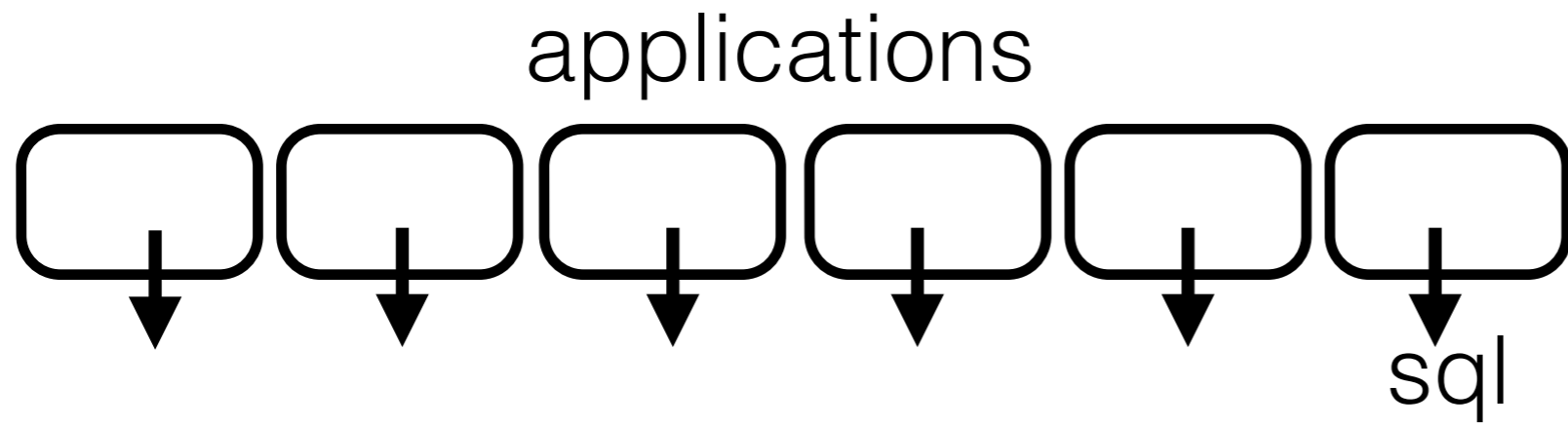
db system

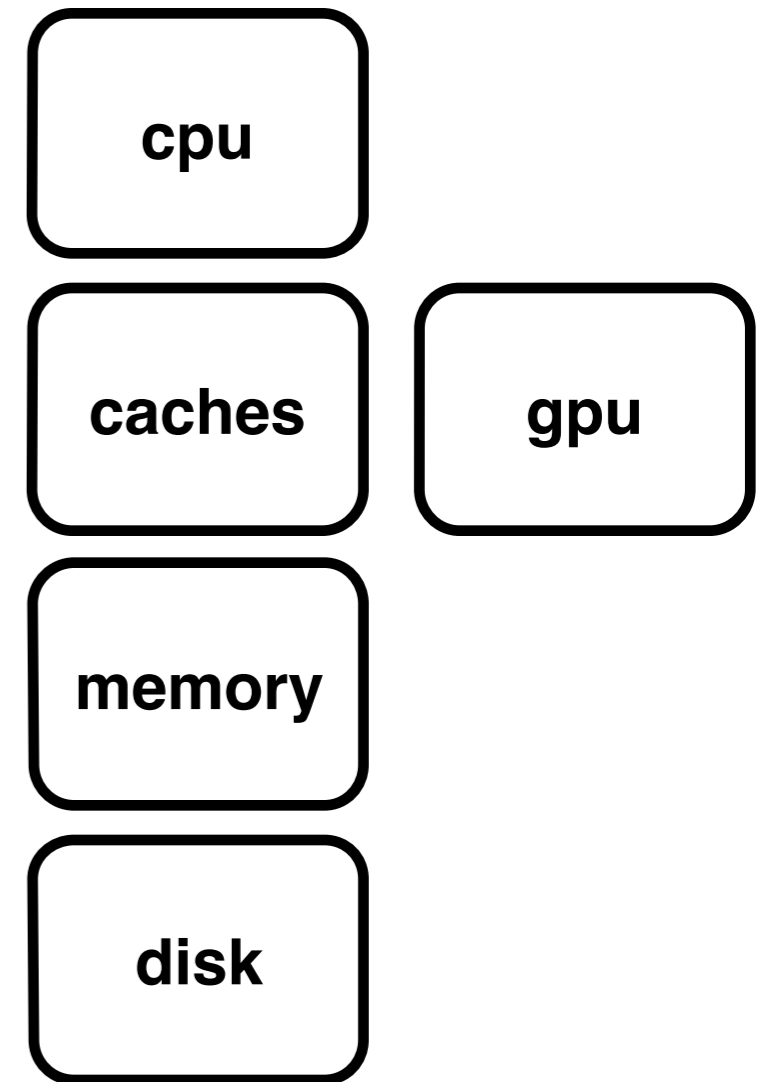
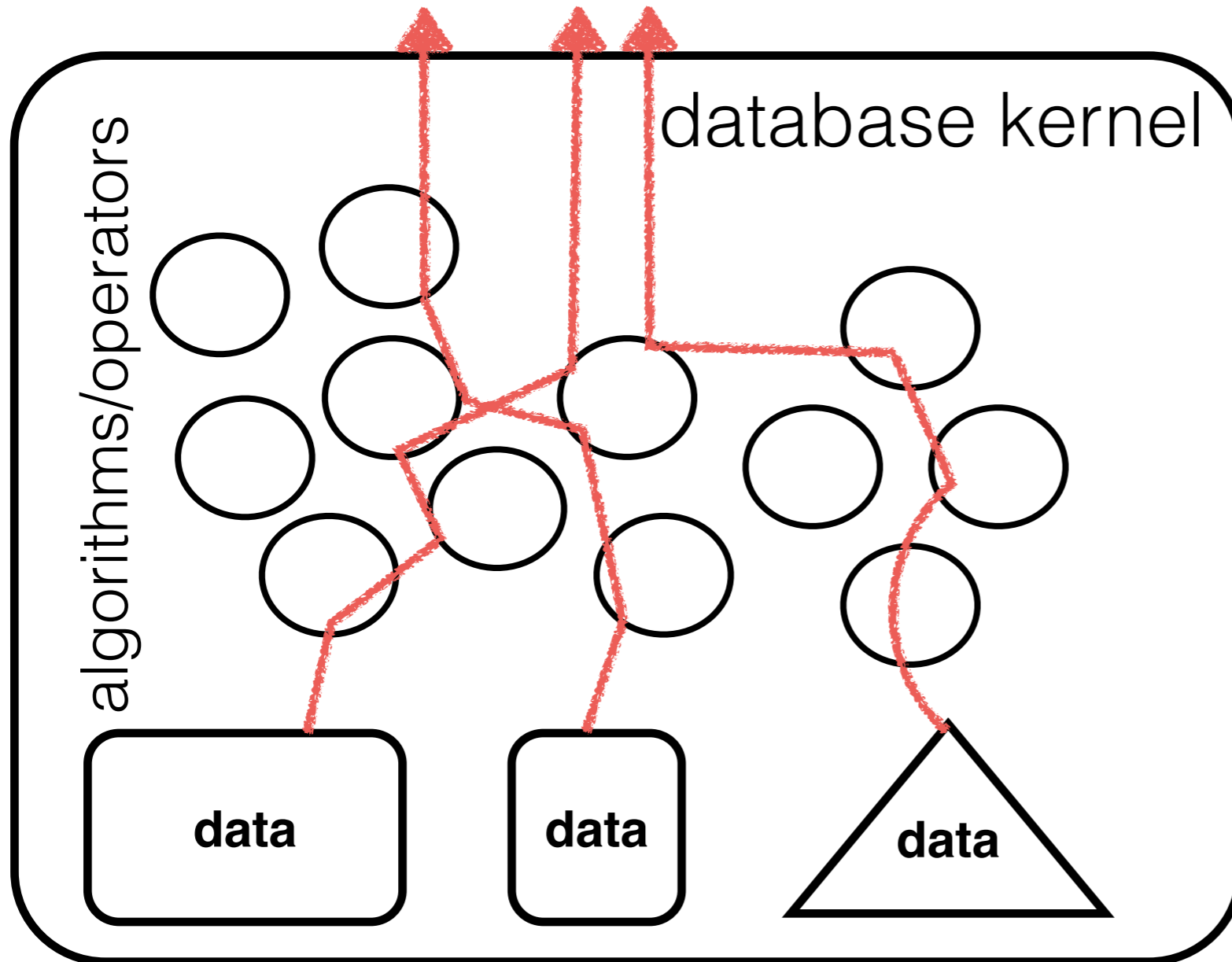
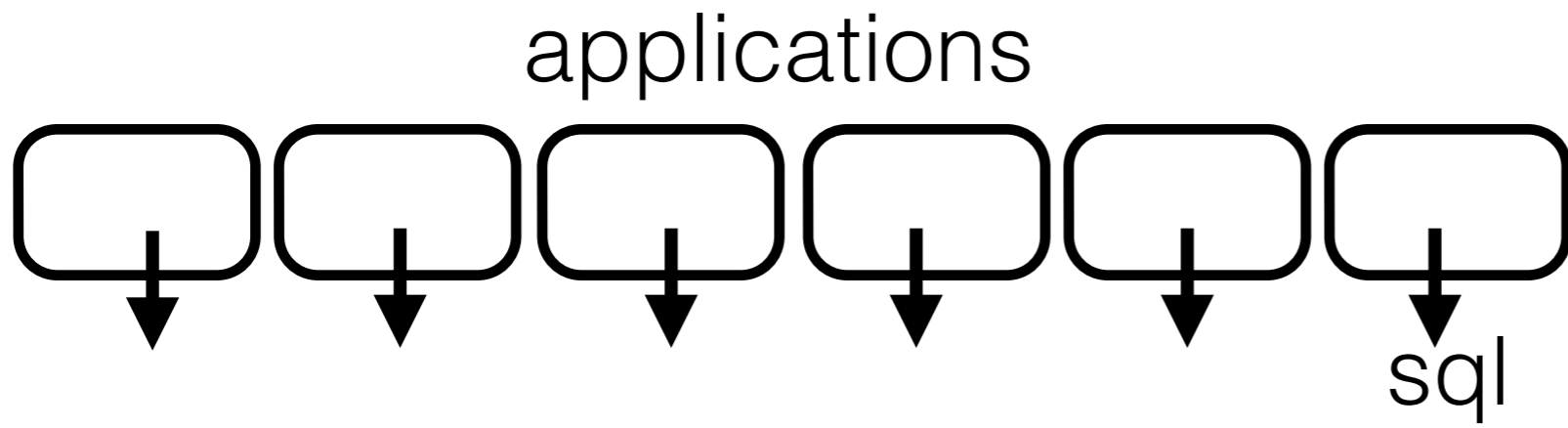
the system decides
"**how**" to best store
and access data

the perfect data system

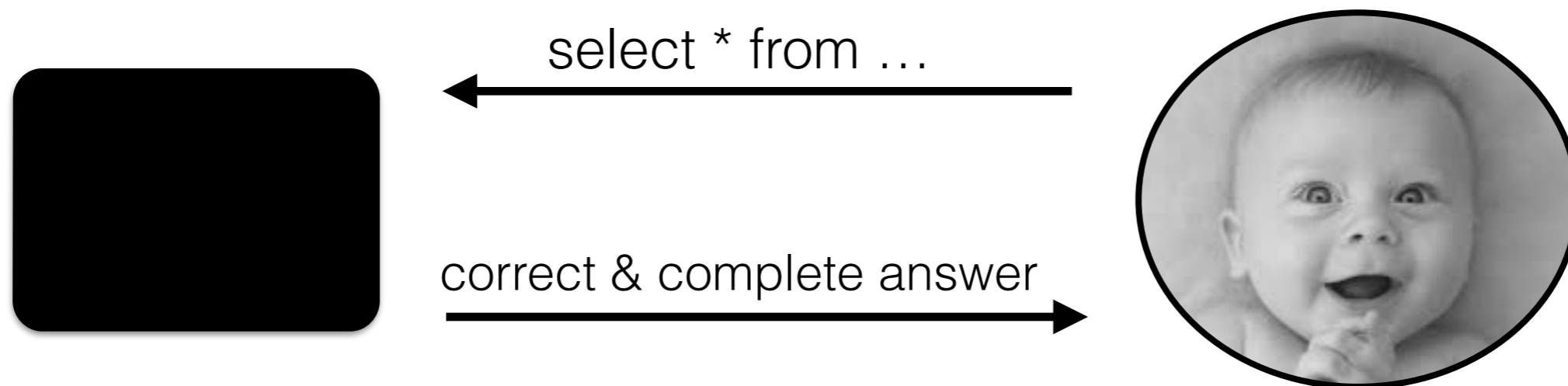




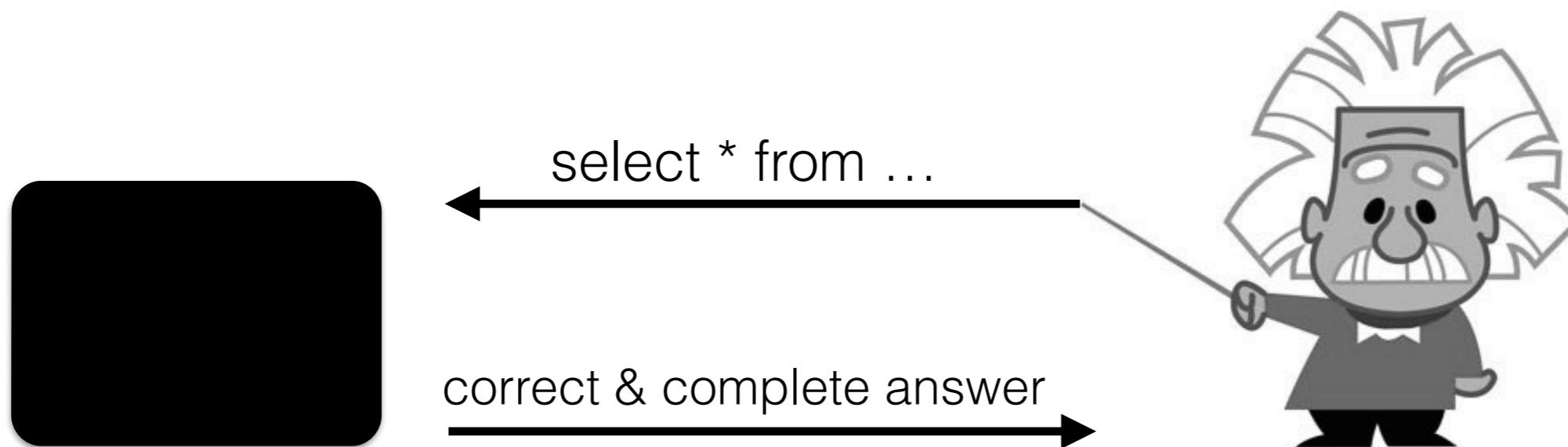




the premise of the black box



the premise of the black box



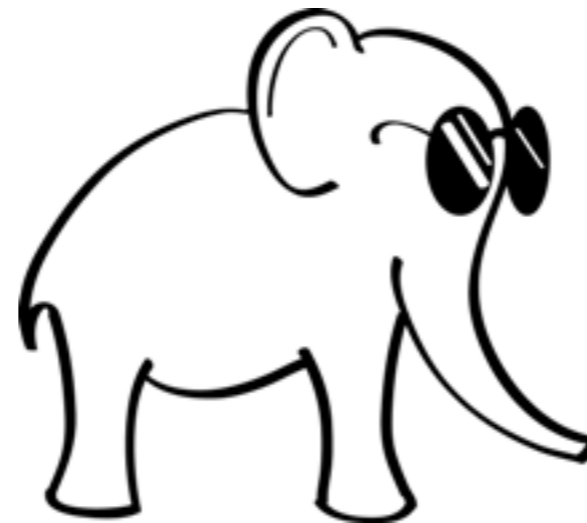
database systems great...

declarative processing, back-end to numerous apps

database systems great...

declarative processing, back-end to numerous apps

but databases have become too heavy and blind!



load

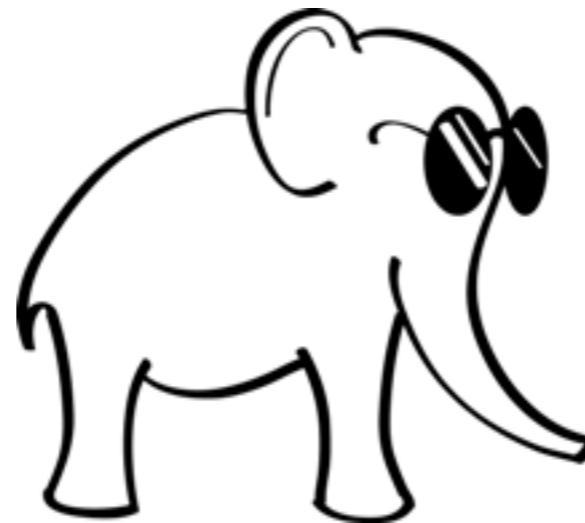
tune

query

timeline

expert users - idle time - workload knowledge

but databases have become too heavy and blind!



load

tune

query

timeline

users/applications
declarative interface
ask what you want



users/applications

need to choose
the proper system

db administrator 1



db administrator 2



data system 1

data system 2

...

usability

tuning

time-to-data

indexing

hardware

size

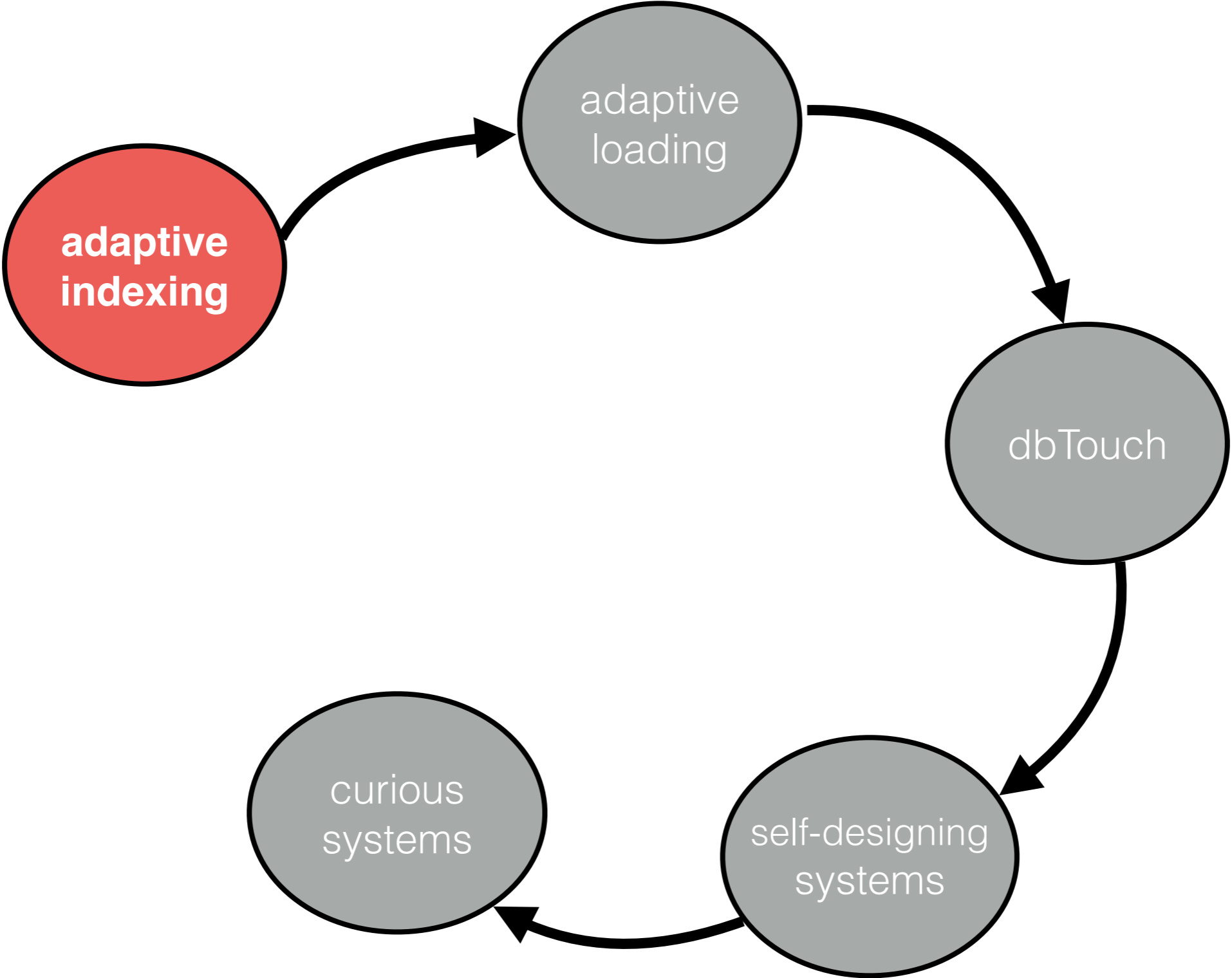
energy

response time

scalability



data systems tailored for data exploration
easy to design - easy to use - fast



indexing



tune= create proper indices offline
performance 10-100X

indexing

load

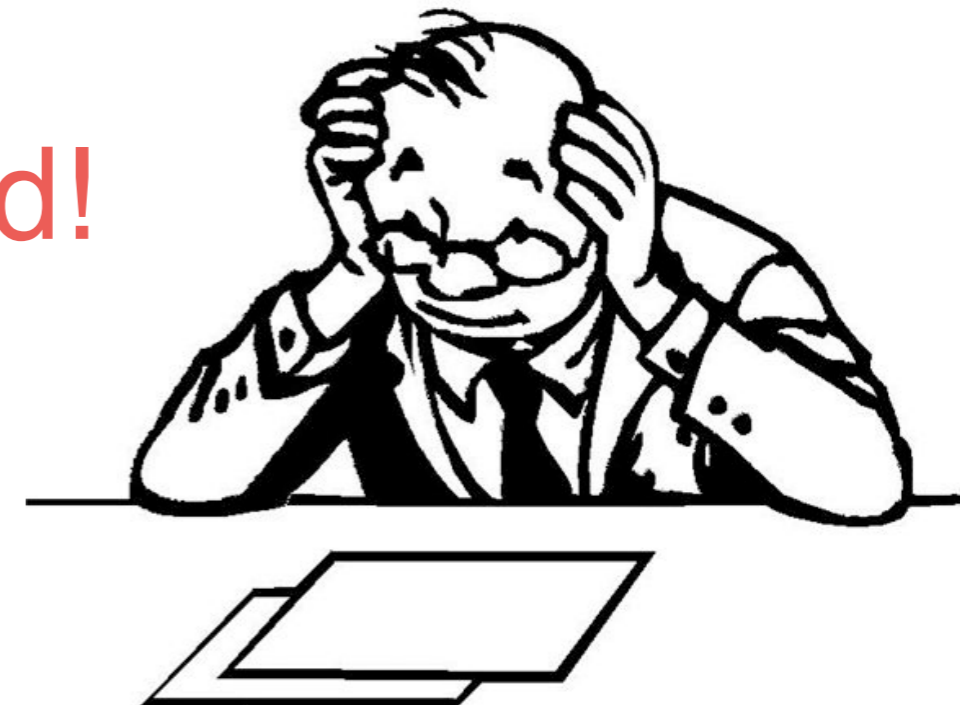
tune

query

tune= create proper indices offline
performance 10-100X

but it depends on workload!

which indices to build?
on which data parts?
and when to build them?



load

tune

query

load

tune

query

timeline

load

tune

query

sample workload

timeline

load

tune

query

sample workload

analyze

timeline →

load

tune

query

sample workload

analyze

create indices

timeline →

load

tune

query

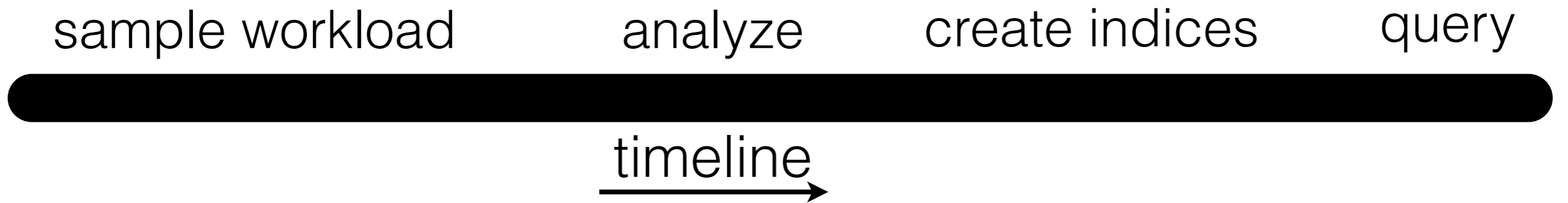
sample workload

analyze

create indices

query

timeline



complex and time consuming process

human administrators + auto-tuning tools

sample workload

analyze

create indices

query

timeline →

complex and time consuming process

big data V's

volume

velocity

variety

veracity

what can go wrong?

not enough space to index all data

not enough idle **time** to finish proper tuning

by the time we finish tuning, the **workload changes**

not enough money - energy - resources

big data V's

volume

velocity

variety

veracity

what can go wrong?

not enough space to index all data

not enough idle **time** to finish proper tuning

by the time we finish tuning, the **workload changes**

not enough money - energy - resources

database cracking

database cracking

~~idle time~~

~~workload
knowledge~~

~~external
tools~~

~~human
control~~

database cracking

auto-tuning database kernels

incremental, adaptive, partial indexing

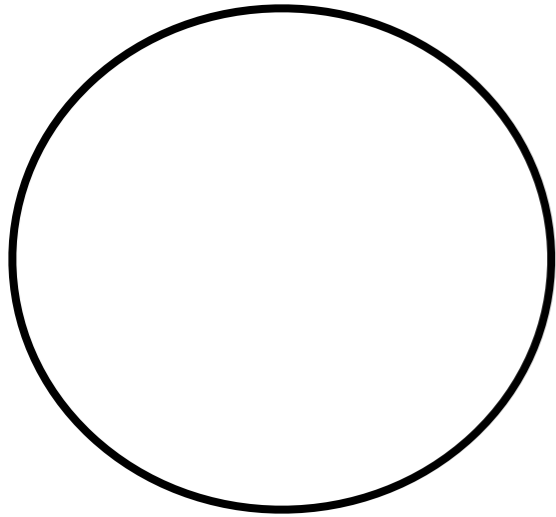
~~idle time~~

~~workload
knowledge~~

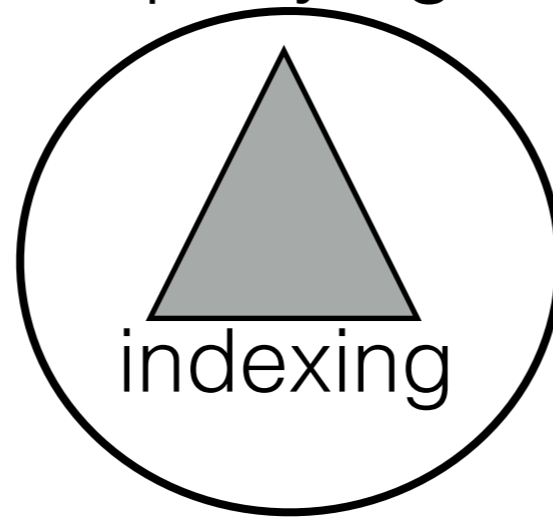
~~external
tools~~

~~human
control~~

initialization



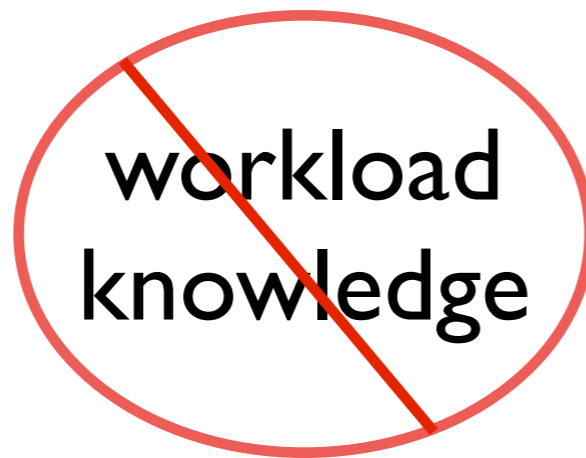
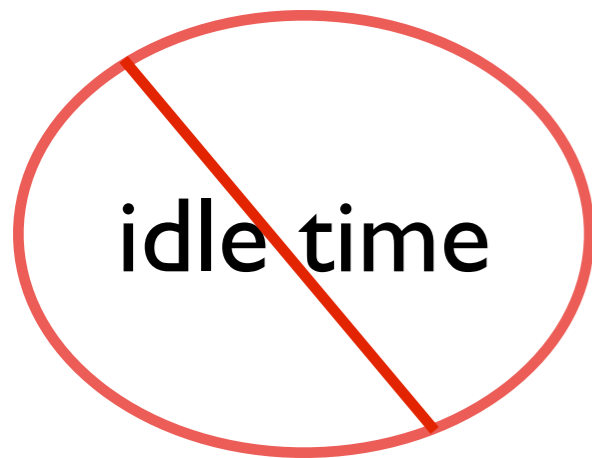
querying



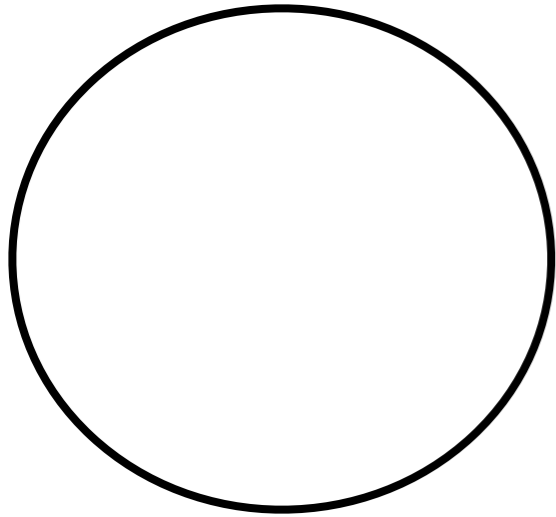
database cracking

auto-tuning database kernels

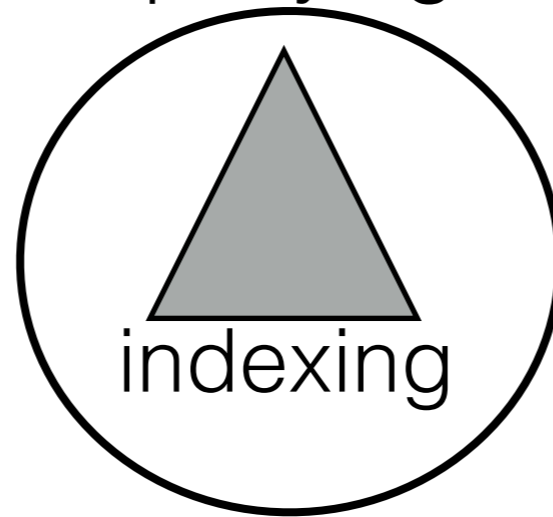
incremental, adaptive, partial indexing



initialization



querying



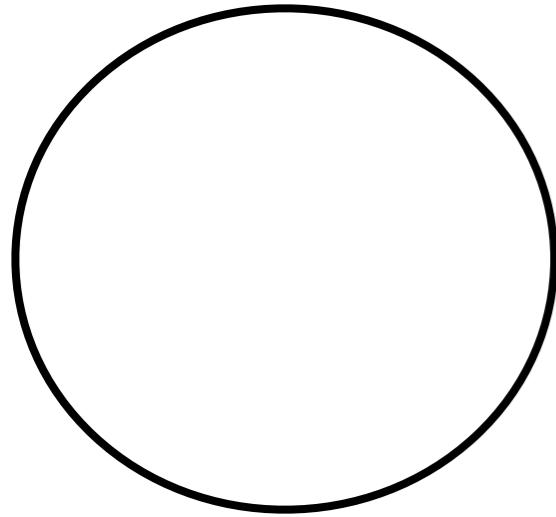
database cracking

auto-tuning database kernels

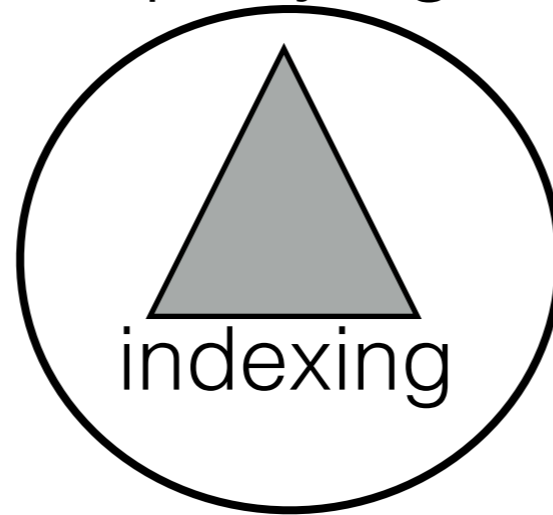
incremental, adaptive, partial indexing



initialization



querying



database cracking

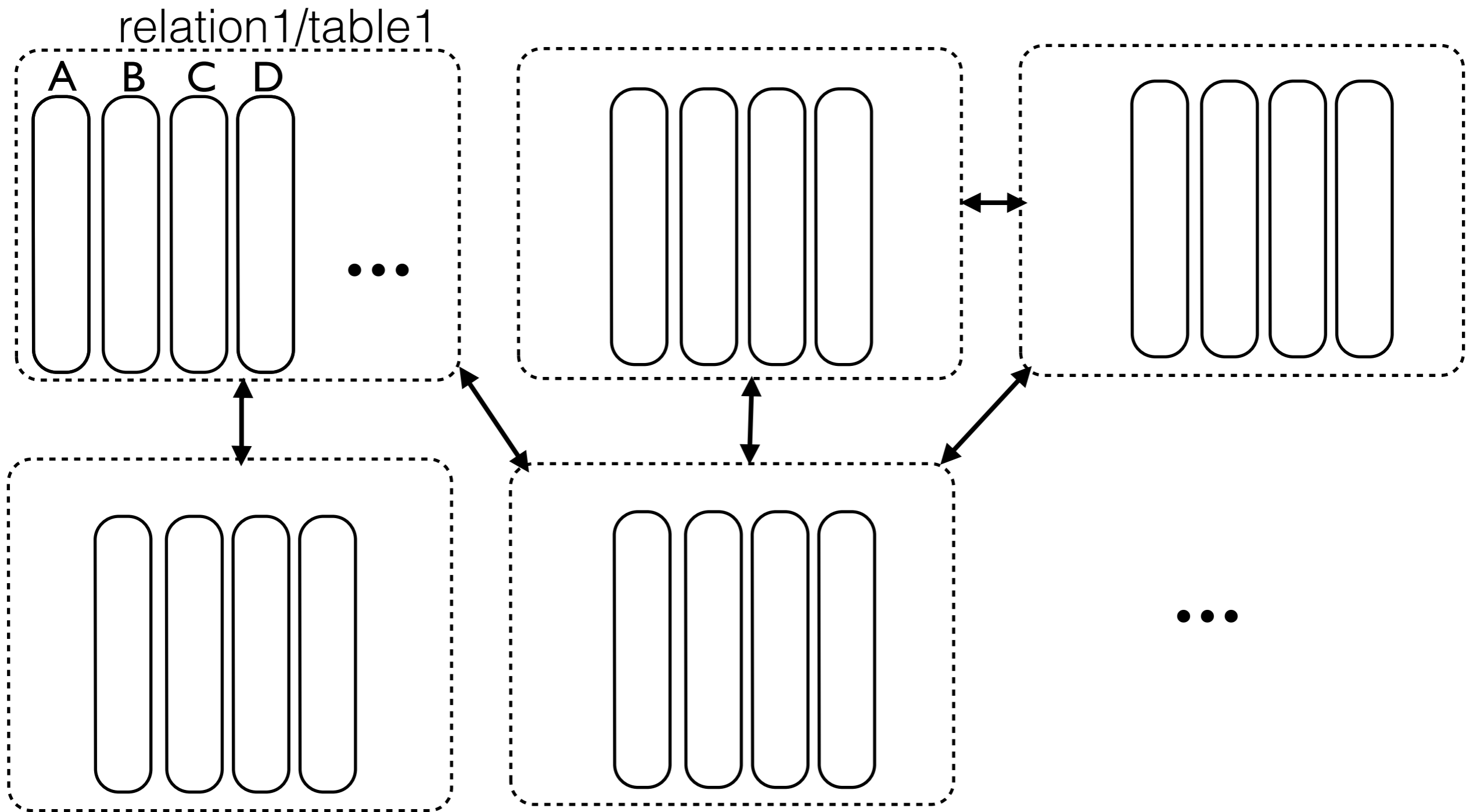
auto-tuning database kernels

incremental, adaptive, partial indexing

**every query is treated as an advice
on how data should be stored**

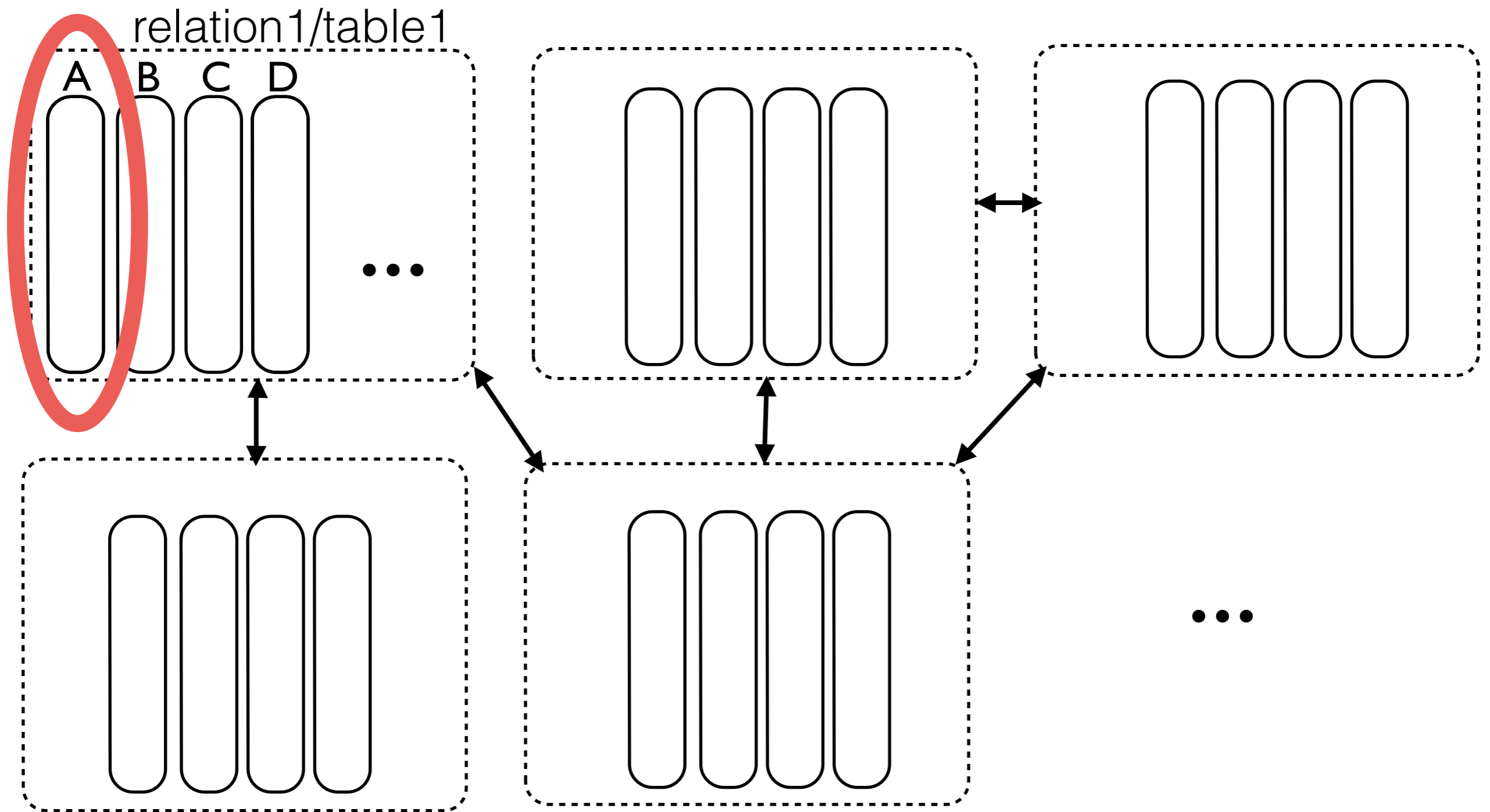
column-store database

a fixed-width and dense array per attribute



column-store database

a fixed-width and dense array per attribute



column A

Q1:
select R.A
from R
where R.A > 10
and R.A < 14

13
16
4
9
2
12
7
1
19
3
14
11
8
6

column A

Q1:
select R.A
from P
where R.A > 10
and R.A < 14

- 13
- 16
- 4
- 9
- 2
- 12
- 7
- 1
- 19
- 3
- 14
- 11
- 8
- 6

column A

Q1:
select R.A
from P
where R.A > 10
and R.A < 14

13
16
4
9
2
12
7
1
19
3
14
11
8
6

sort

1
2
3
4
6
7
8
9
11
12
13
14
16
19

Q1:
select R.A
from P
where R.A > 10
and R.A < 14

column A

- 13
- 16
- 4
- 9
- 2
- 12
- 7
- 1
- 19
- 3
- 14
- 11
- 8
- 6

sort
→

binary
search

- 1
- 2
- 3
- 4
- 6
- 7
- 8
- 9
- 11
- 12
- 13
- 14
- 16
- 19

Q1:
select R.A
from P
where R.A > 10
and R.A < 14

column A

- 13
- 16
- 4
- 9
- 2
- 12
- 7
- 1
- 19
- 3
- 14
- 11
- 8
- 6

sort
→

binary
search

- 1
 - 2
 - 3
 - 4
 - 6
 - 7
 - 8
 - 9
 - 11
 - 12
 - 13
 - 14
 - 16
 - 19
- result

Q1:
select R.A
from P
where R.A > 10
and R.A < 14

column A

- 13
- 16
- 4
- 9
- 2
- 12
- 7
- 1
- 19
- 3
- 14
- 11
- 8
- 6

sort
→

binary
search

- 1
- 2
- 3
- 4
- 6
- 7
- 8
- 9
- 11
- 12
- 13
- 14
- 16
- 19

time
+
knowledge

result

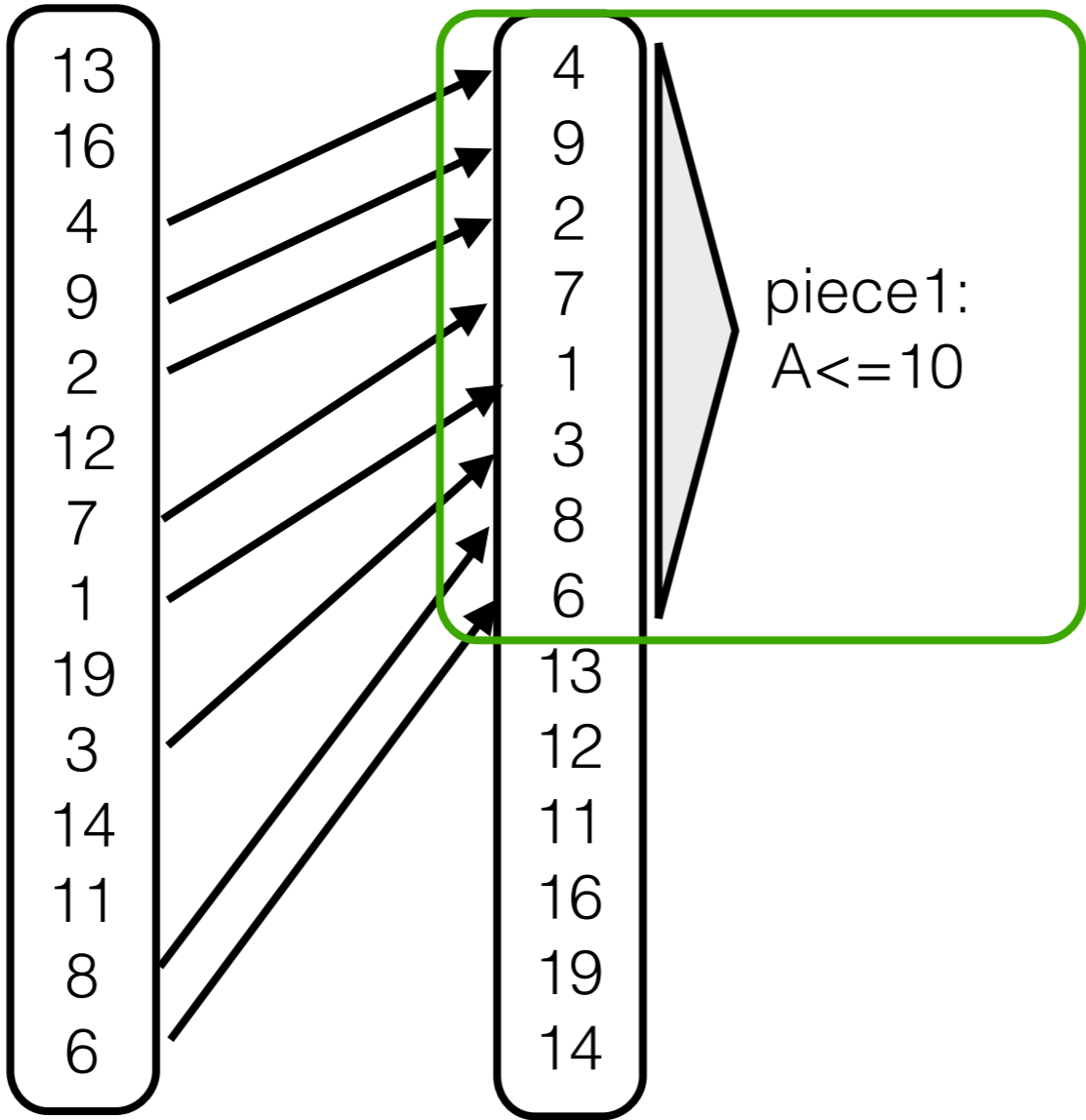
column A

Q1:
select R.A
from R
where R.A > 10
and R.A < 14

13
16
4
9
2
12
7
1
19
3
14
11
8
6

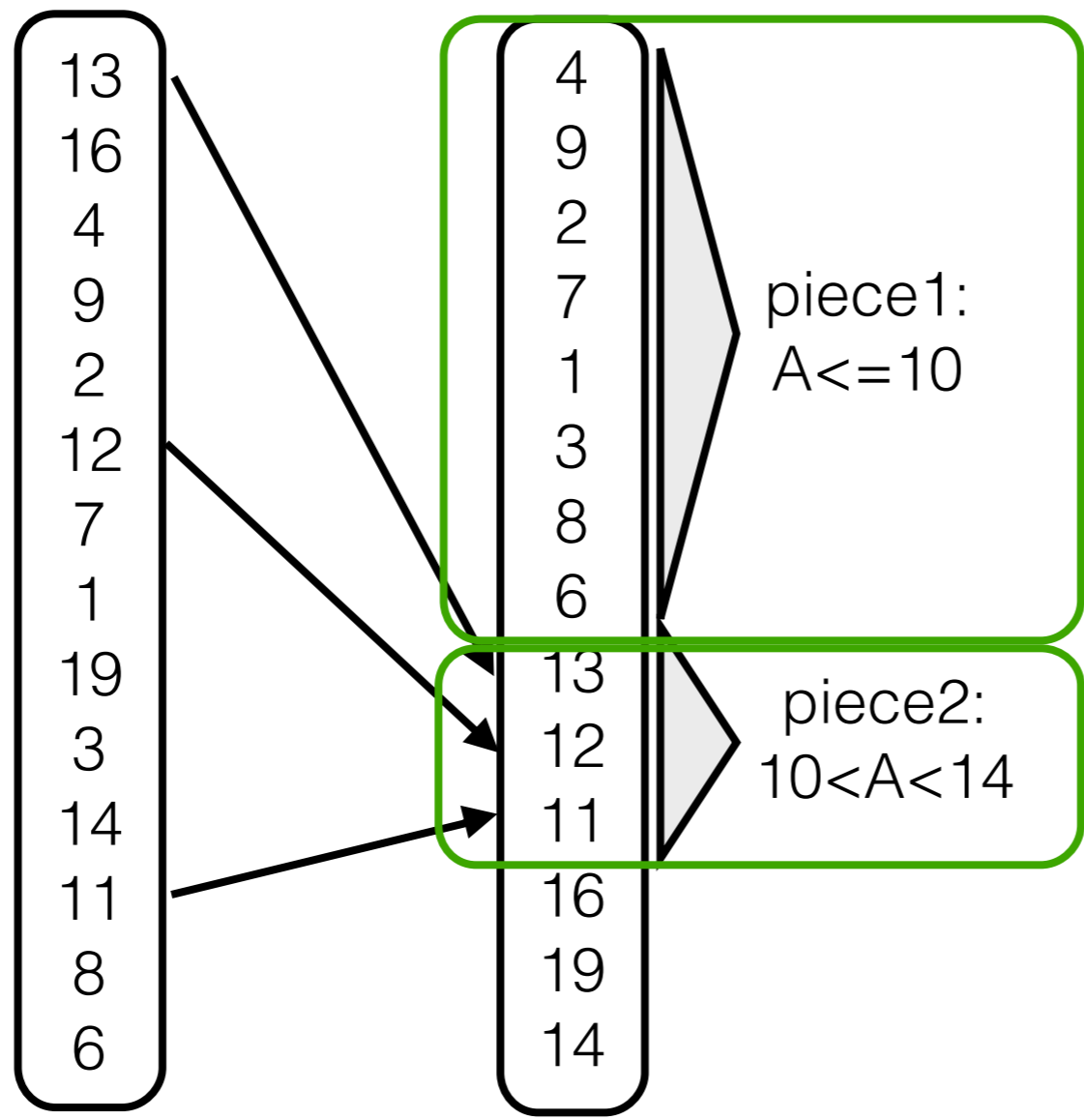
column A

Q1:
select R.A
from R
where R.A > 10
and R.A < 14



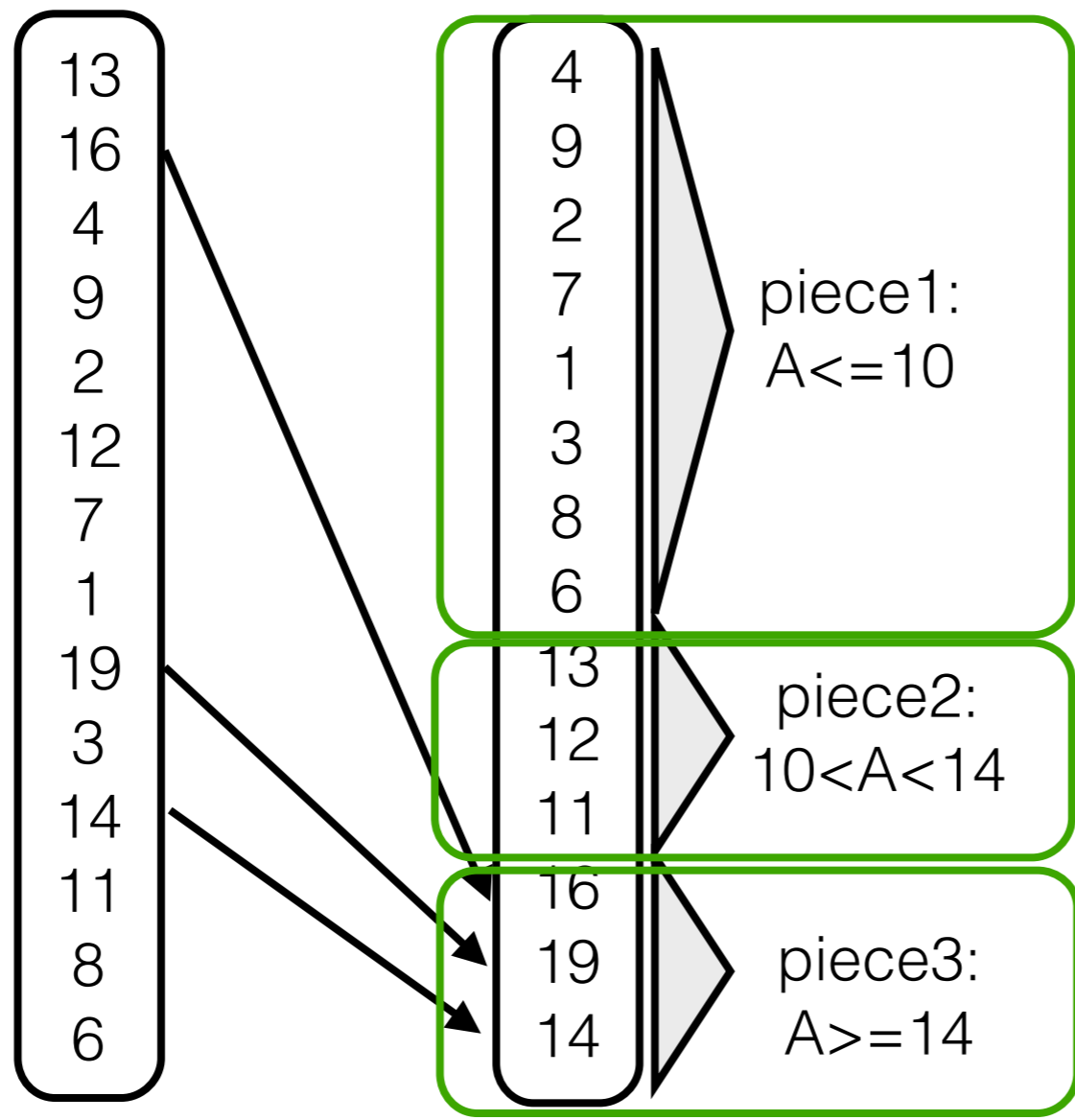
Q1:
select R.A
from R
where R.A > 10
and R.A < 14

column A



Q1:
select R.A
from R
where R.A > 10
and R.A < 14

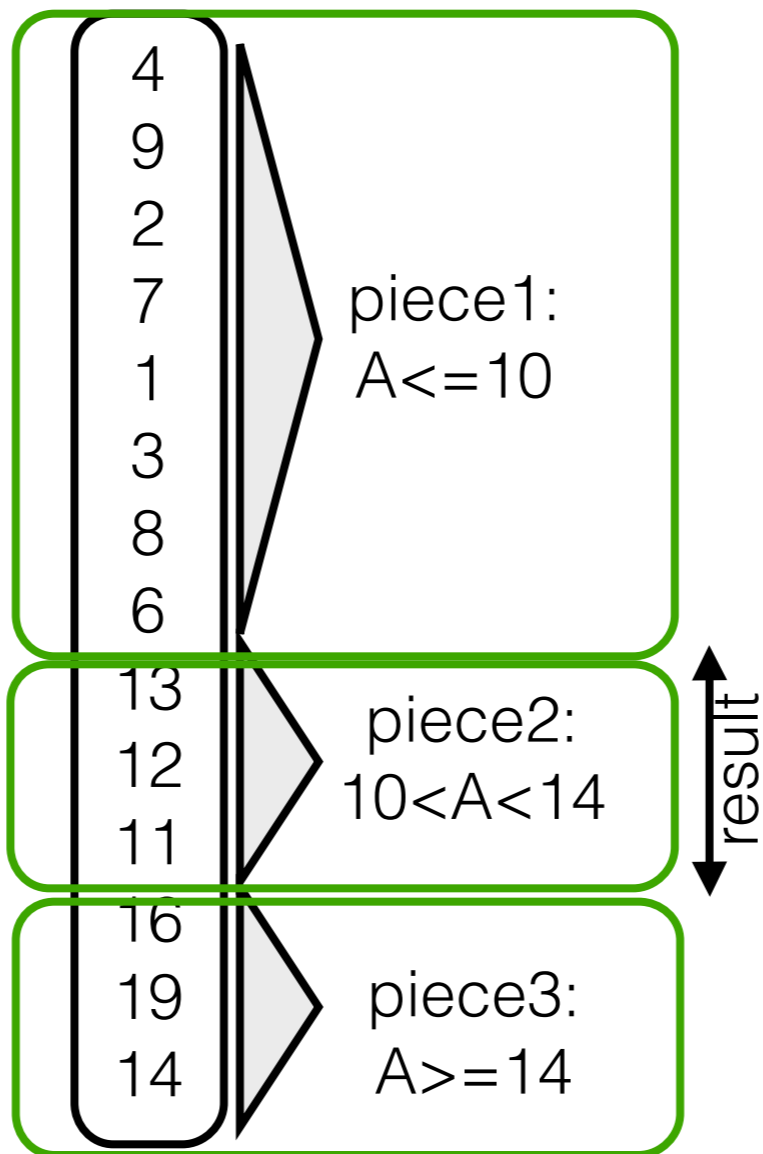
column A



Q1:
select R.A
from R
where R.A > 10
and R.A < 14

column A

13
16
4
9
2
12
7
1
19
3
14
11
8
6

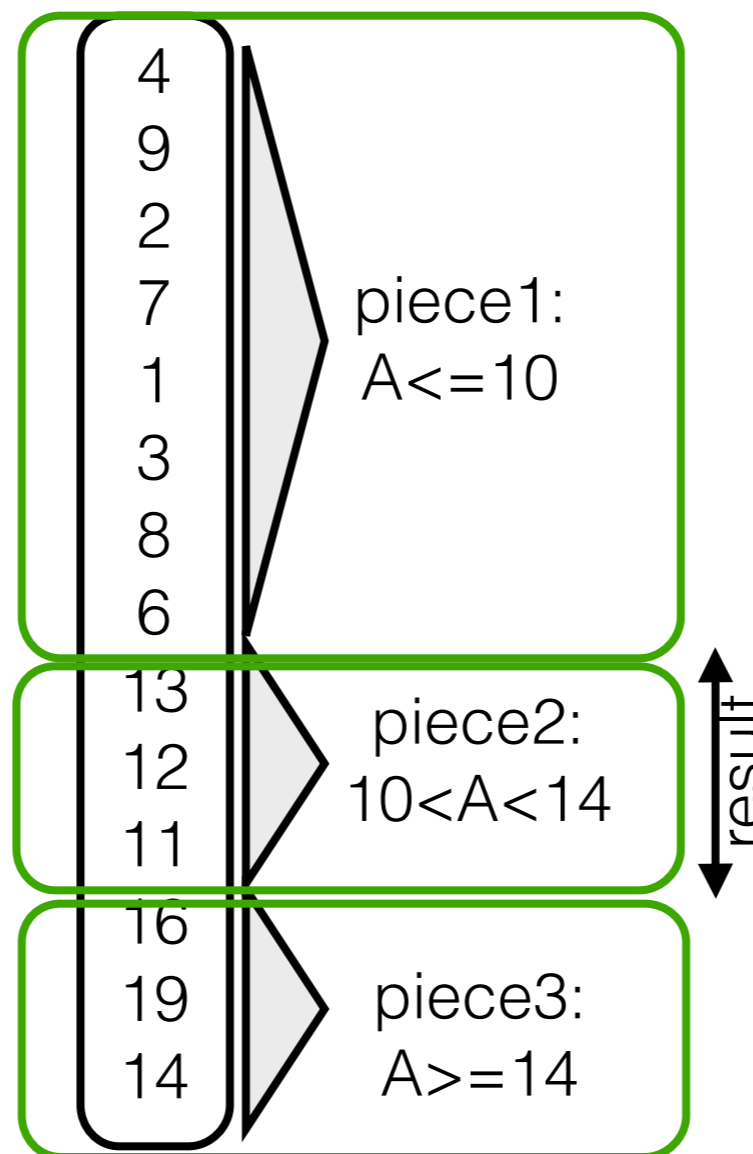


gain knowledge on how data is organized

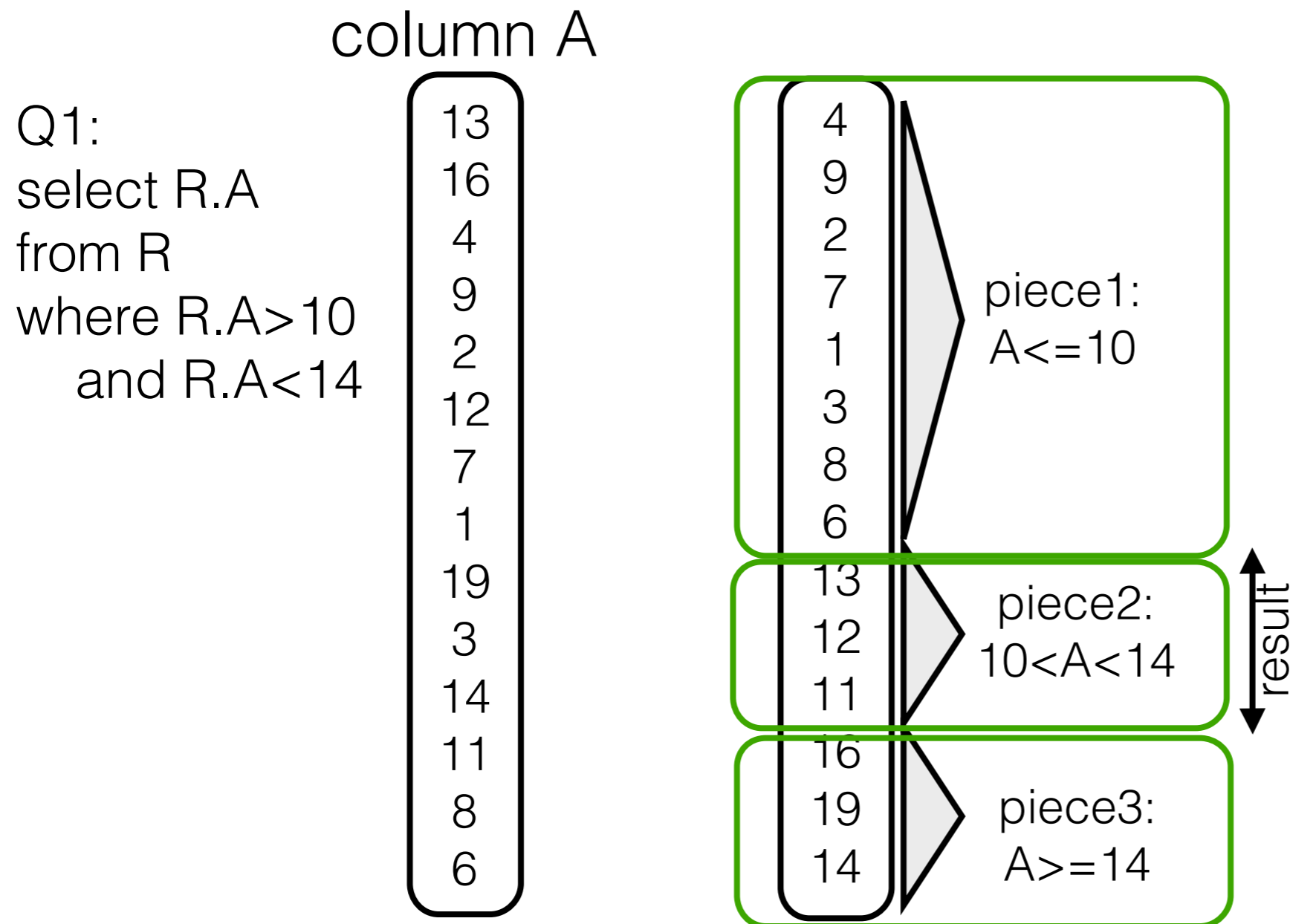
column A

Q1:
select R.A
from R
where R.A > 10
and R.A < 14

13
16
4
9
2
12
7
1
19
3
14
11
8
6



gain knowledge on how data is organized



dynamically/on-the-fly within the select-operator

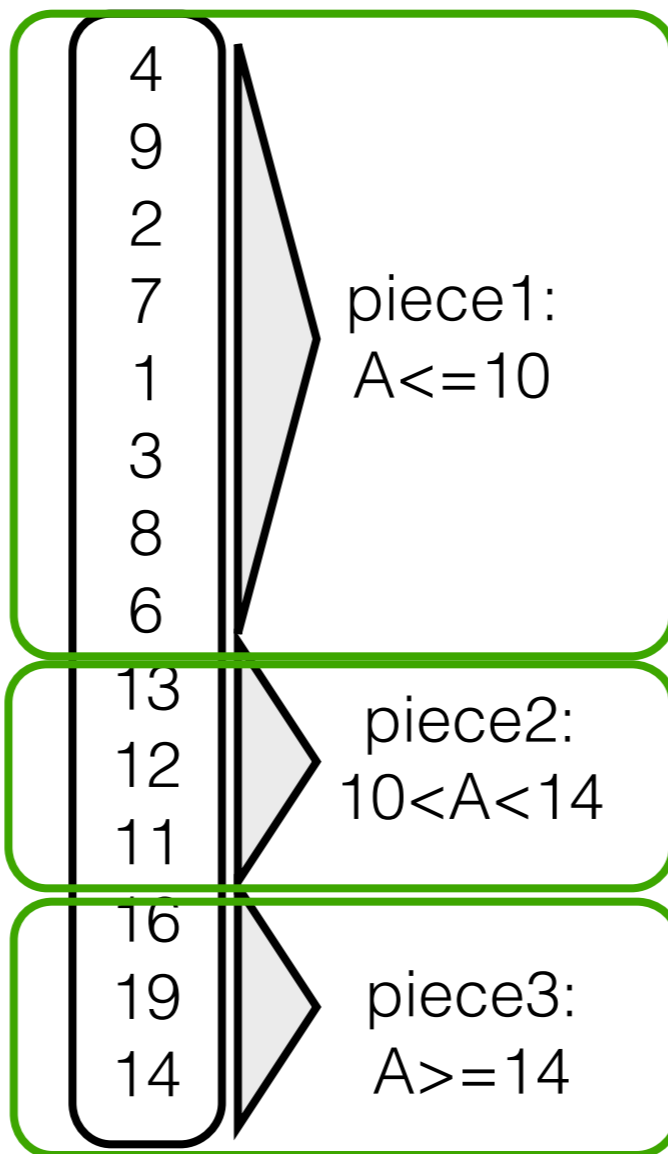
Database Cracking CIDR 2007

column A

Q1:
select R.A
from R
where R.A > 10
and R.A < 14

Q2:
select R.A
from R
where R.A > 7
and R.A <= 16

13
16
4
9
2
12
7
1
19
3
14
11
8
6



dynamically/on-the-fly within the select-operator

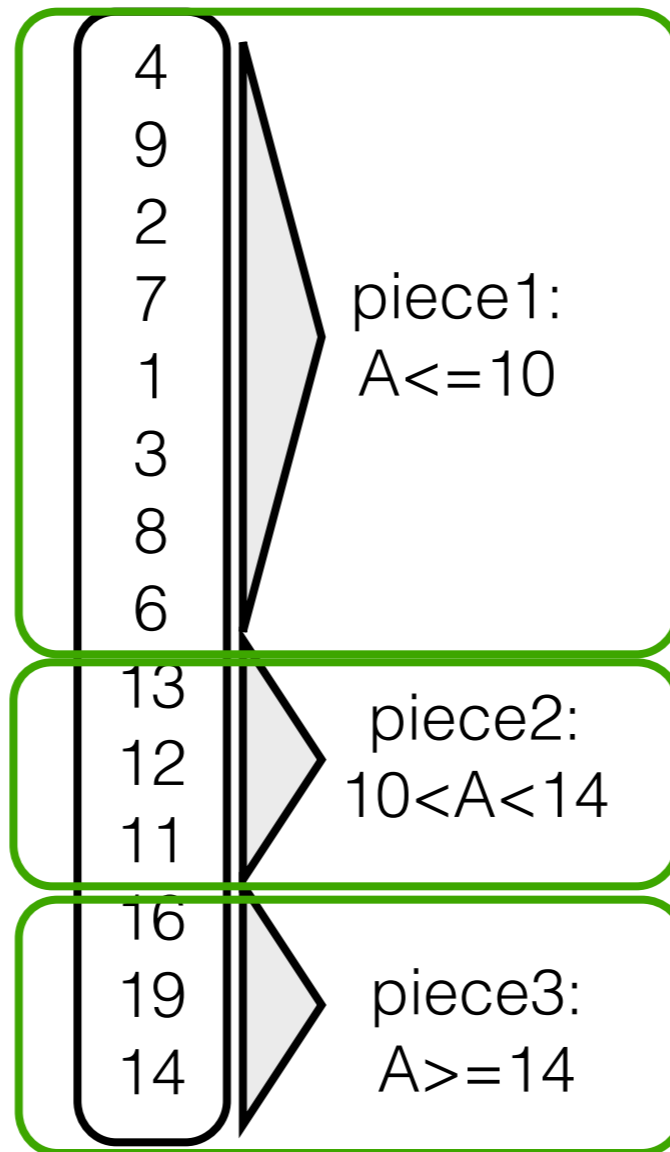
Database Cracking CIDR 2007

column A

Q1:
select R.A
from R
where R.A > 10
and R.A < 14

Q2:
select R.A
from R
where R.A > 7
and R.A <= 16

13
16
4
9
2
12
7
1
19
3
14
11
8
6



dynamically/on-the-fly within the select-operator

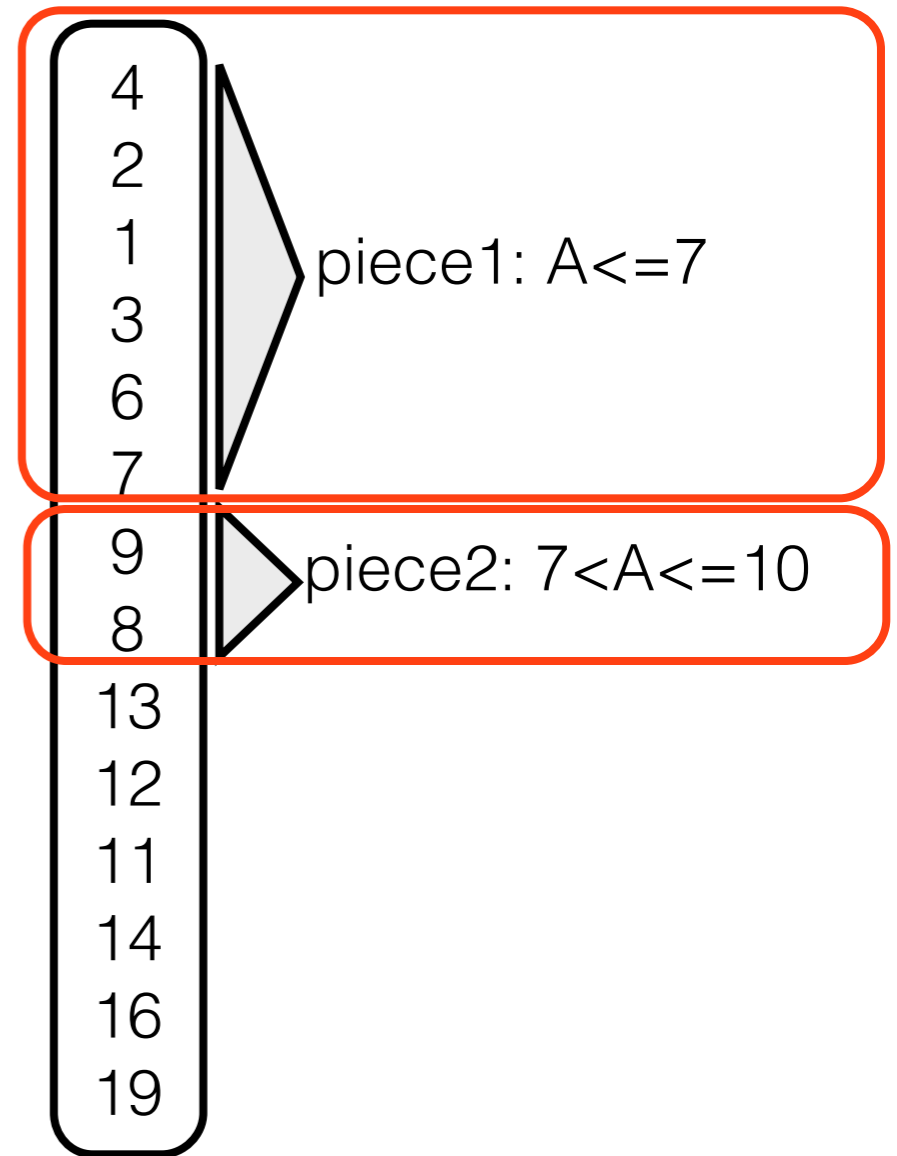
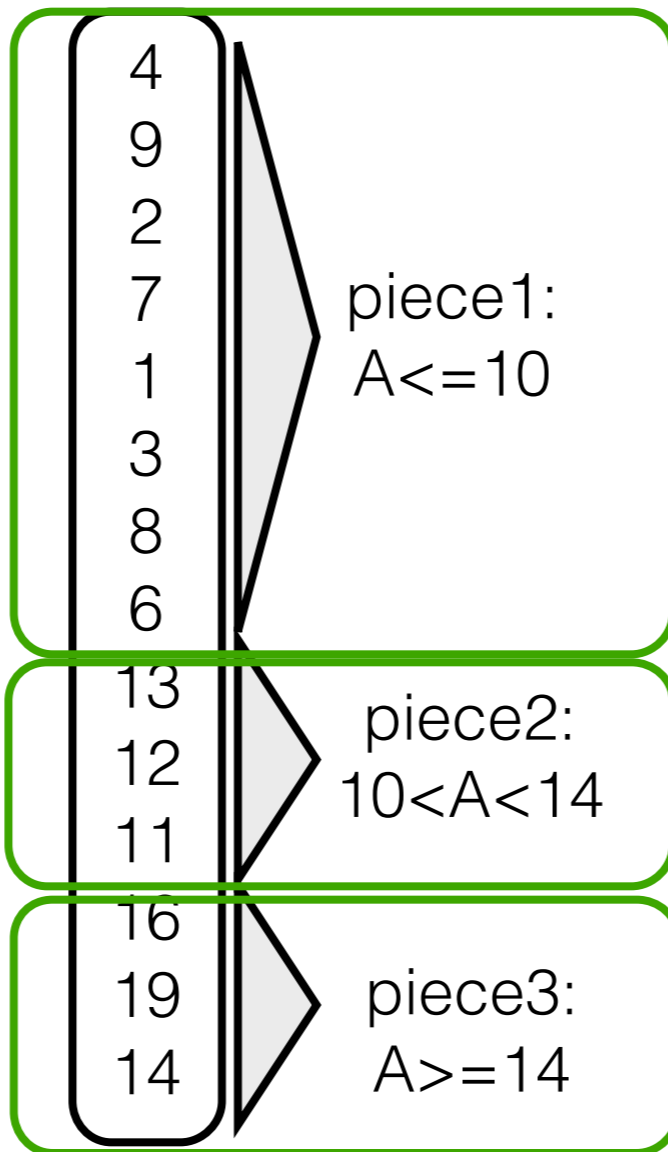
Database Cracking CIDR 2007

column A

Q1:
select R.A
from R
where R.A > 10
and R.A < 14

Q2:
select R.A
from R
where R.A > 7
and R.A <= 16

13
16
4
9
2
12
7
1
19
3
14
11
8
6



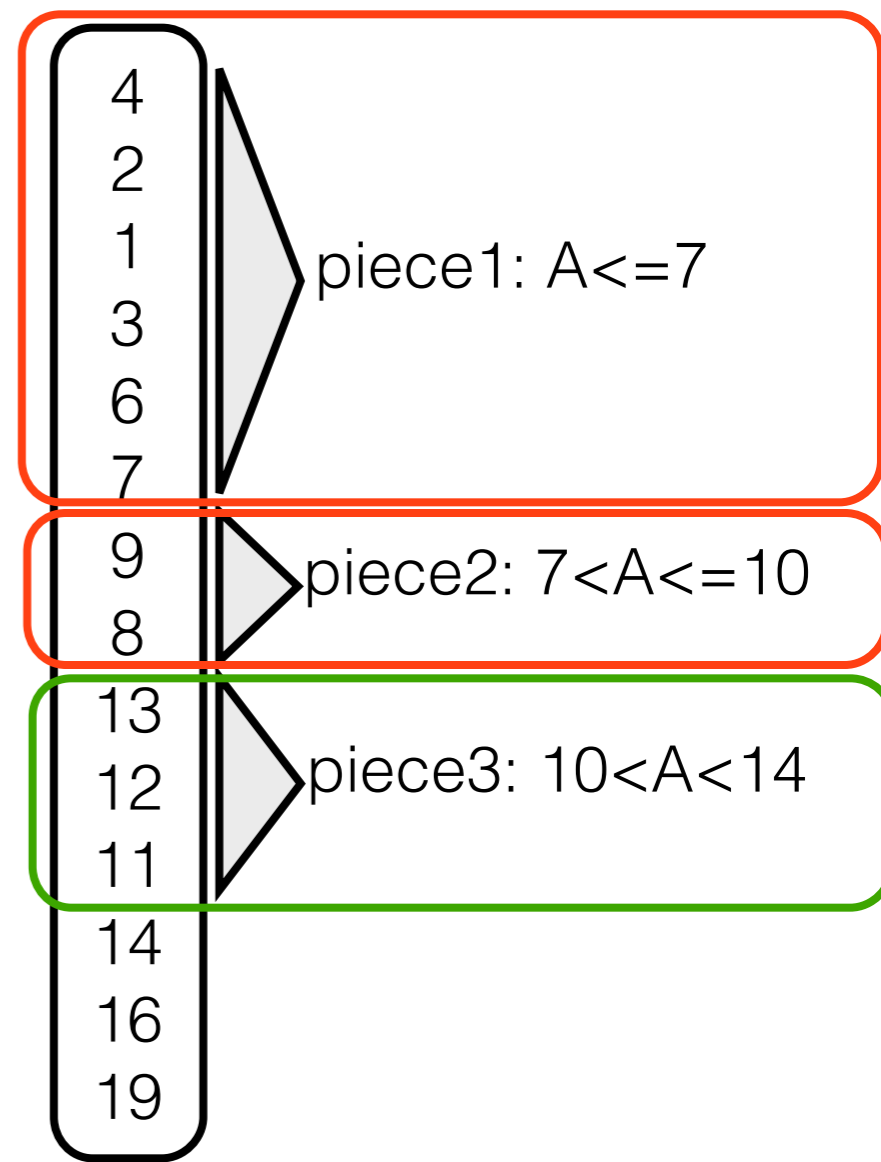
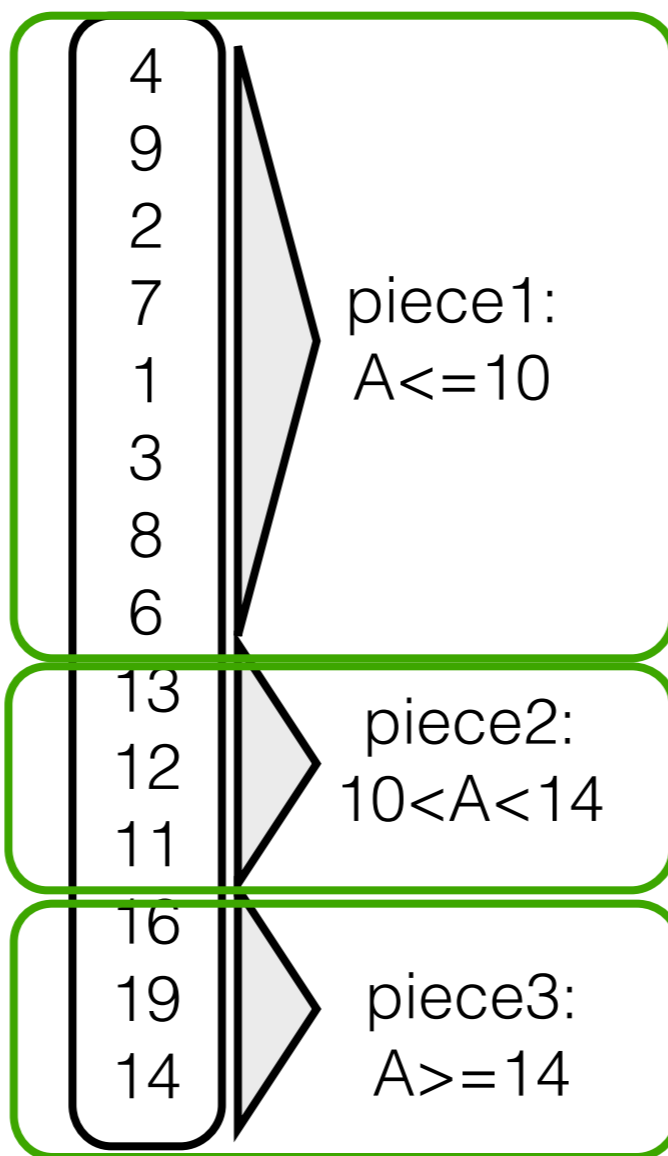
dynamically/on-the-fly within the select-operator

column A

Q1:
select R.A
from R
where R.A > 10
and R.A < 14

Q2:
select R.A
from R
where R.A > 7
and R.A <= 16

13
16
4
9
2
12
7
1
19
3
14
11
8
6



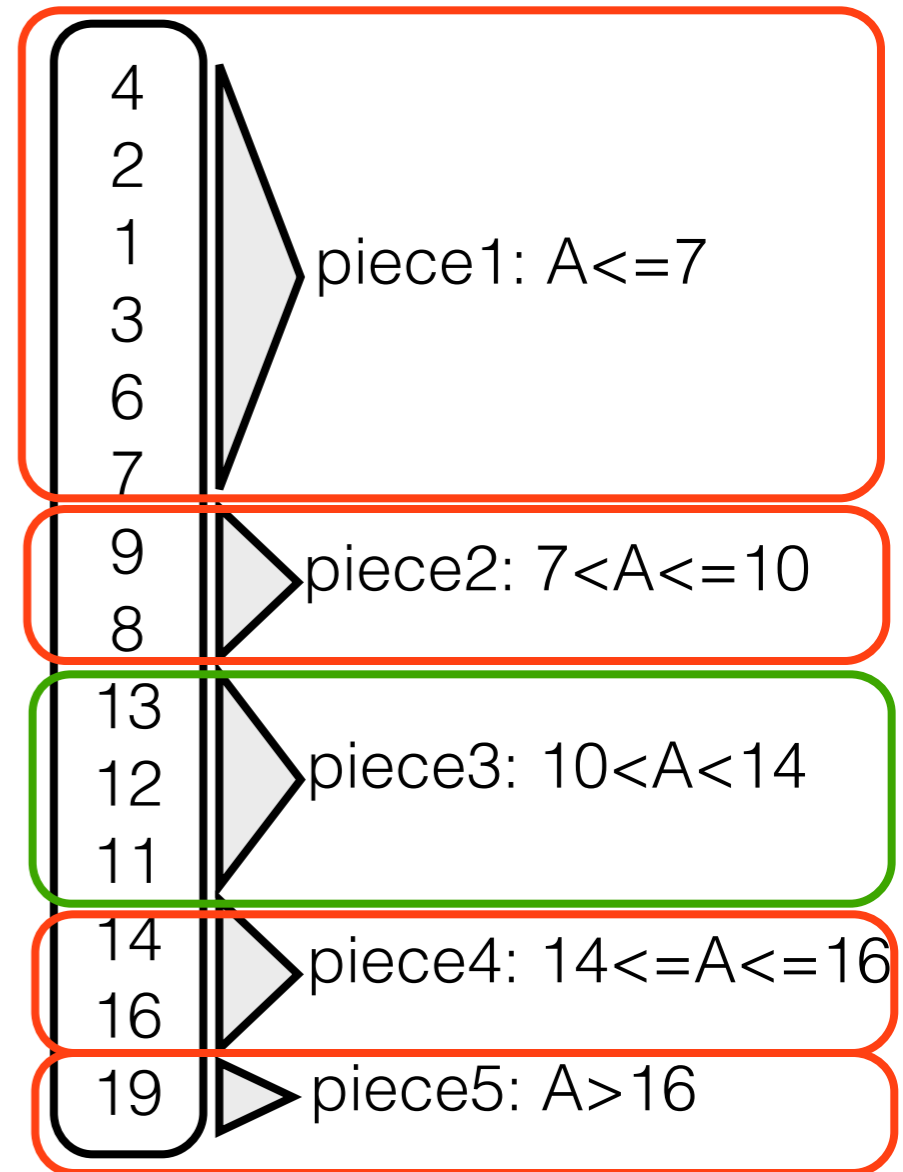
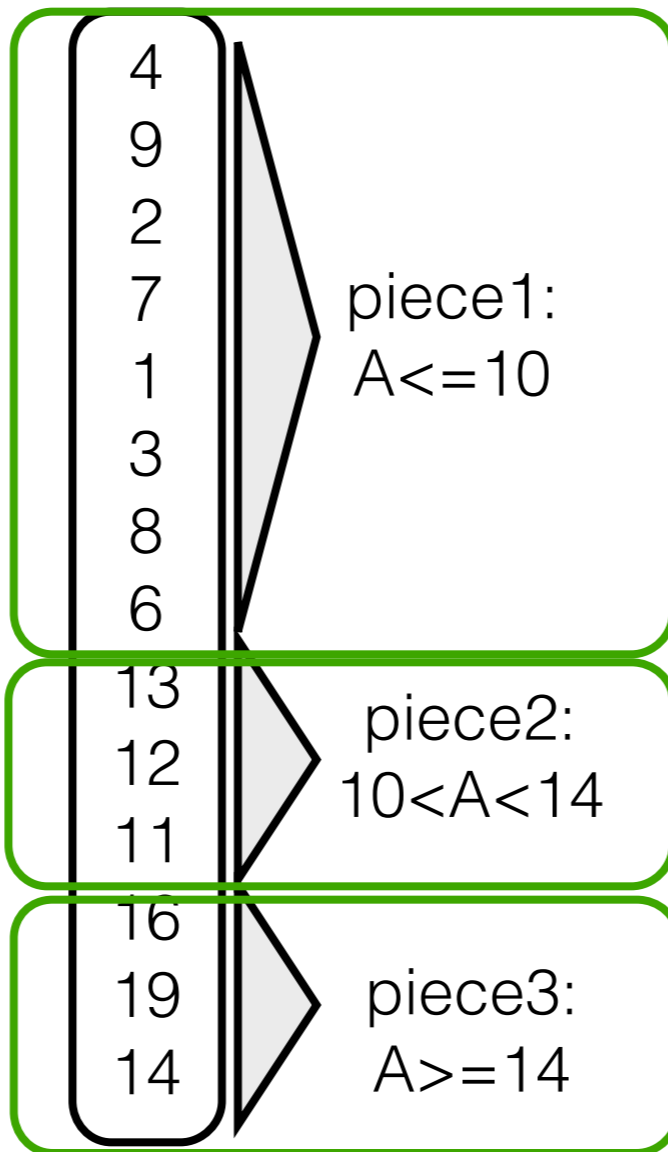
dynamically/on-the-fly within the select-operator

column A

Q1:
select R.A
from R
where R.A > 10
and R.A < 14

Q2:
select R.A
from R
where R.A > 7
and R.A <= 16

13
16
4
9
2
12
7
1
19
3
14
11
8
6



dynamically/on-the-fly within the select-operator

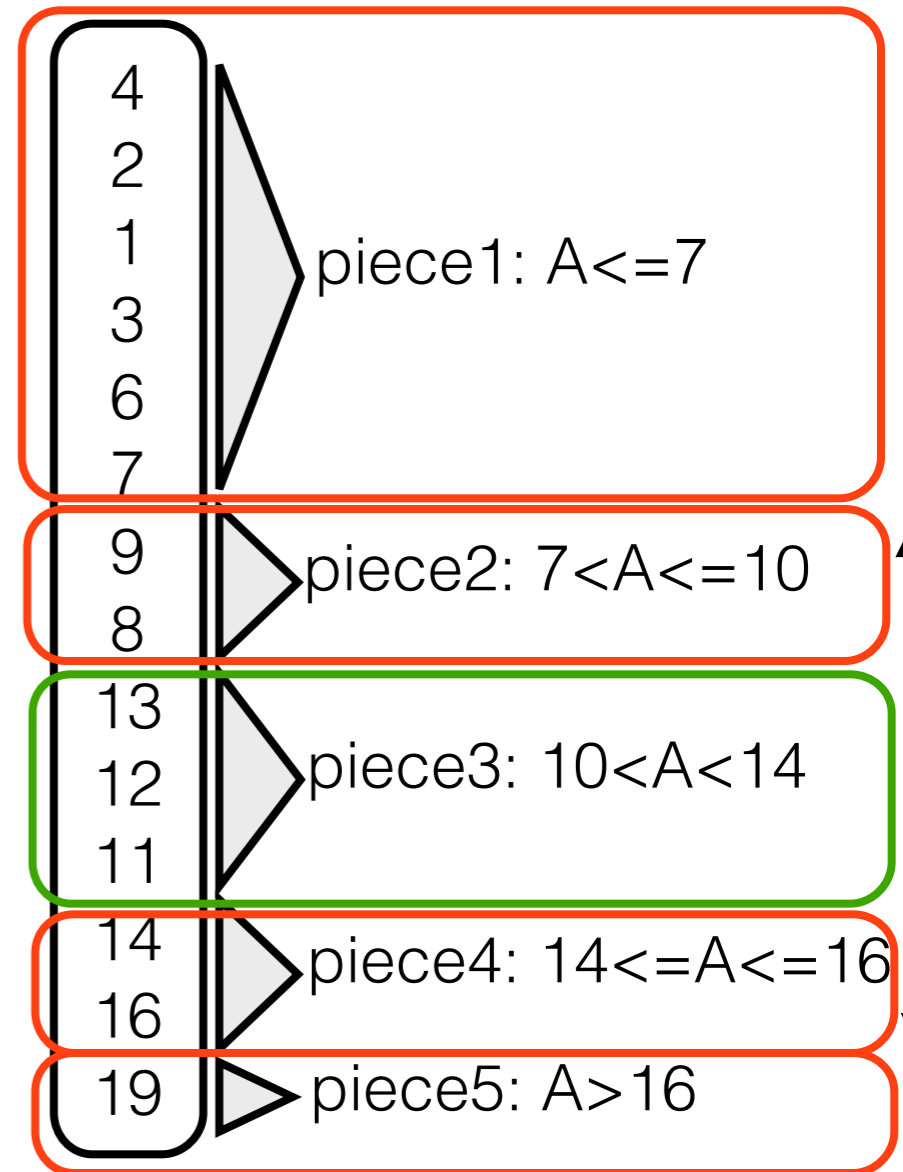
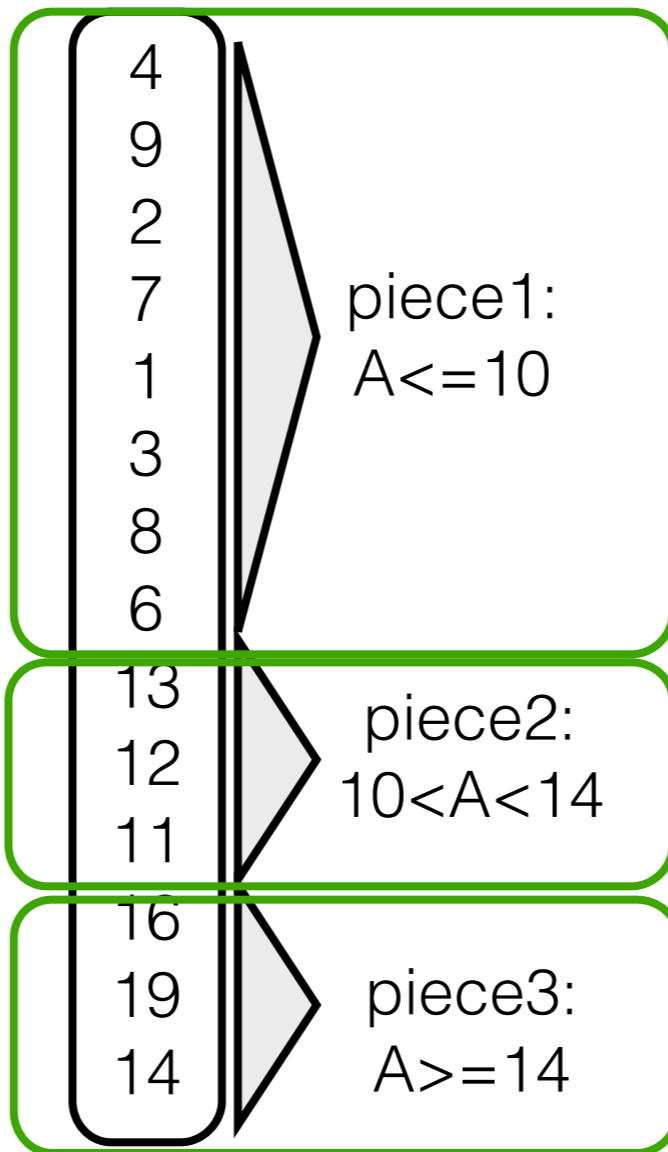
Database Cracking CIDR 2007

column A

Q1:
select R.A
from R
where R.A > 10
and R.A < 14

Q2:
select R.A
from R
where R.A > 7
and R.A <= 16

13
16
4
9
2
12
7
1
19
3
14
11
8
6



dynamically/on-the-fly within the select-operator

Database Cracking CIDR 2007

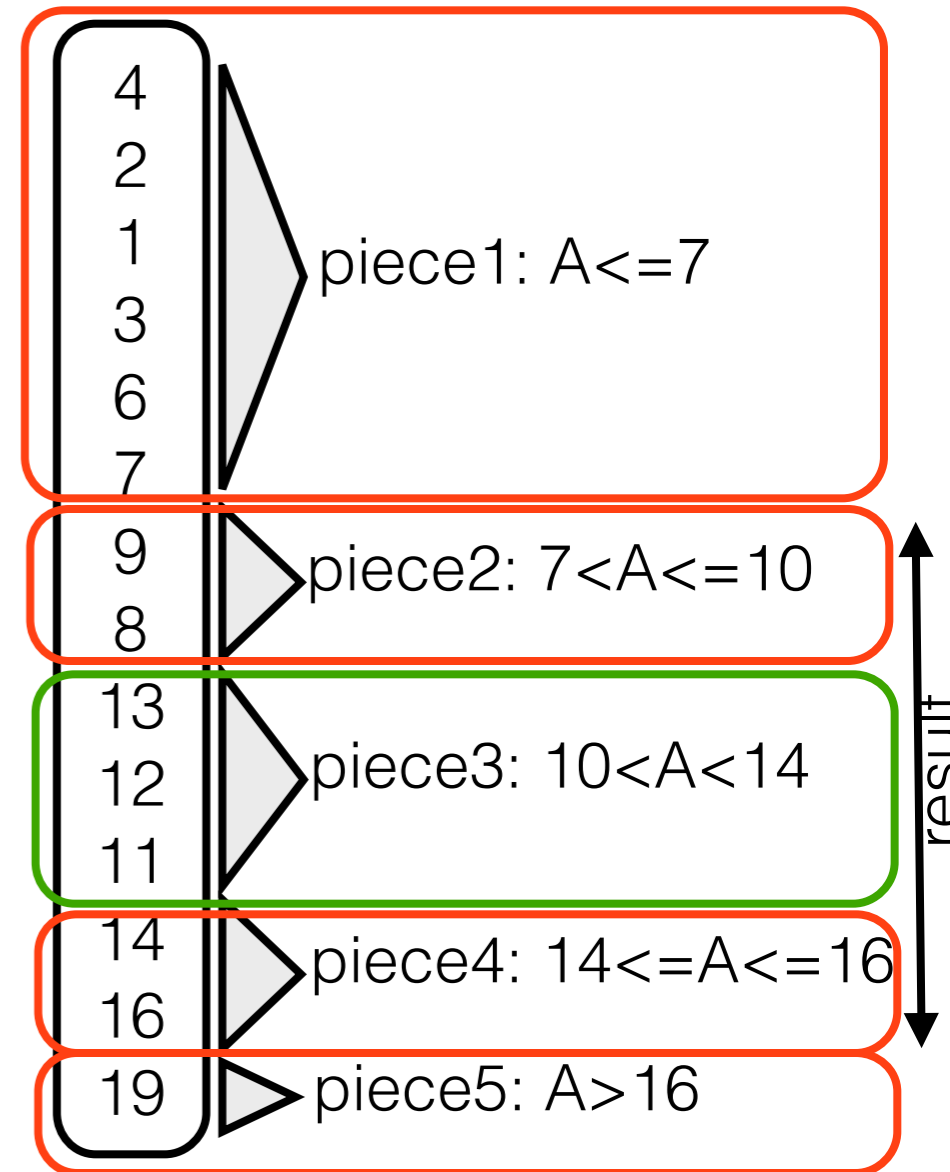
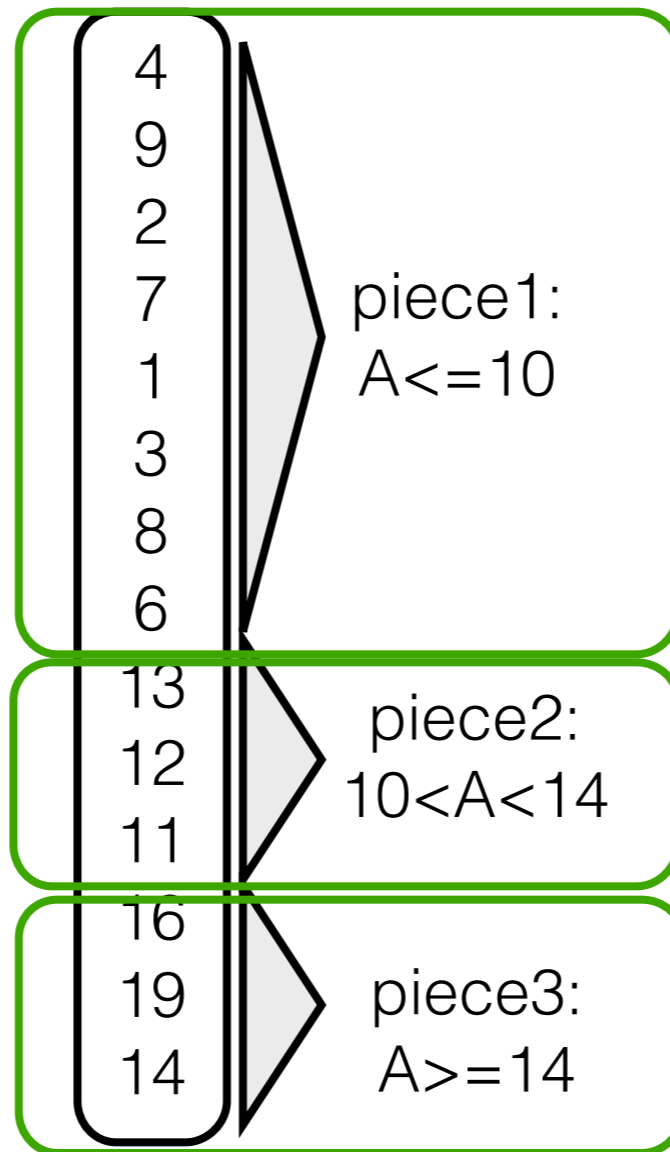
the more we crack, the more we learn

column A

Q1:
select R.A
from R
where R.A > 10
and R.A < 14

Q2:
select R.A
from R
where R.A > 7
and R.A <= 16

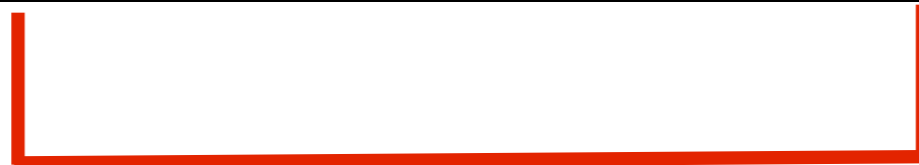
13
16
4
9
2
12
7
1
19
3
14
11
8
6



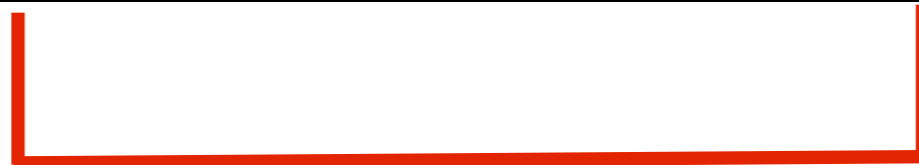
dynamically/on-the-fly within the select-operator

Database Cracking CIDR 2007

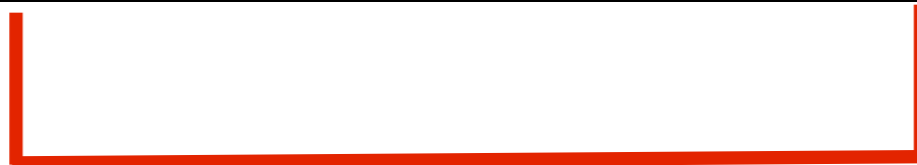




select [15,55]



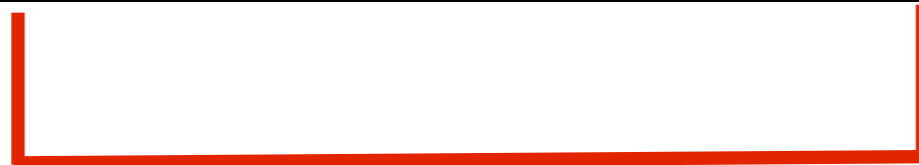
select [15,55]



select [15,55]

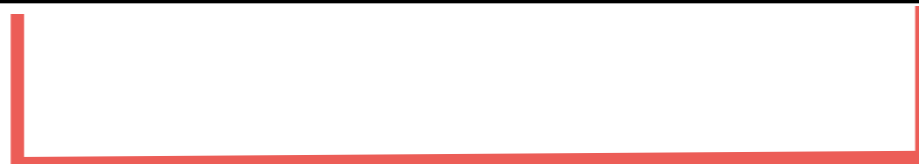
10 20 30 40 50 60



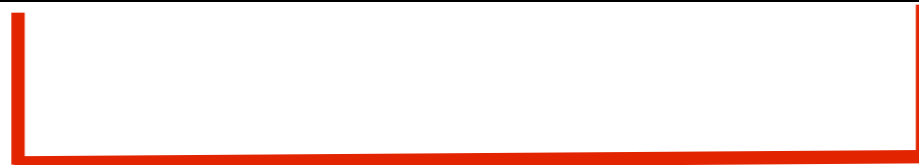


select [15,55]

10 20 30 40 50 60

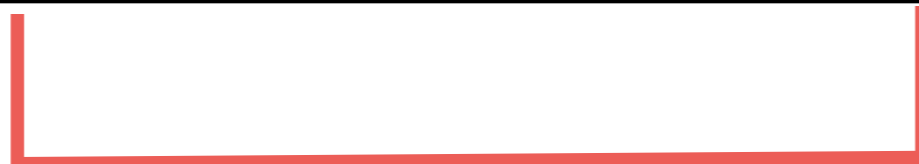


select [15,55]



select [15,55]

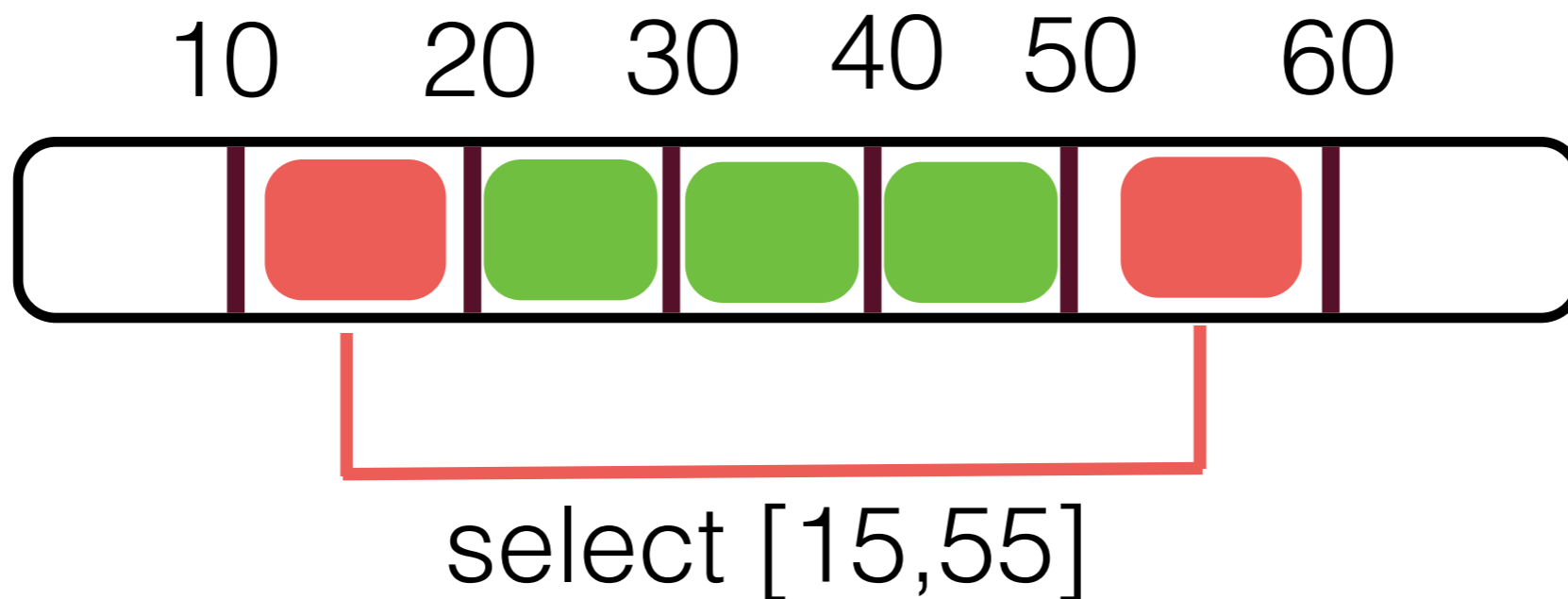
10 20 30 40 50 60




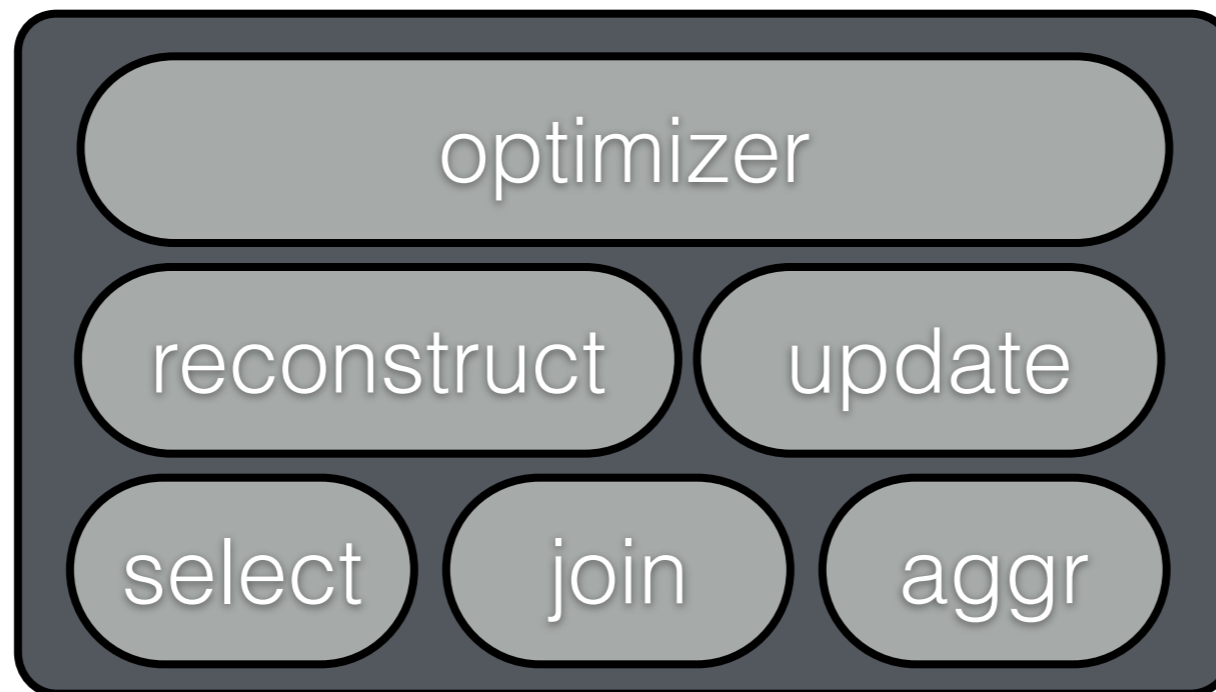
select [15,55]

touch at most two pieces at a time

pieces become smaller and smaller



implemented in 
open-source column-store

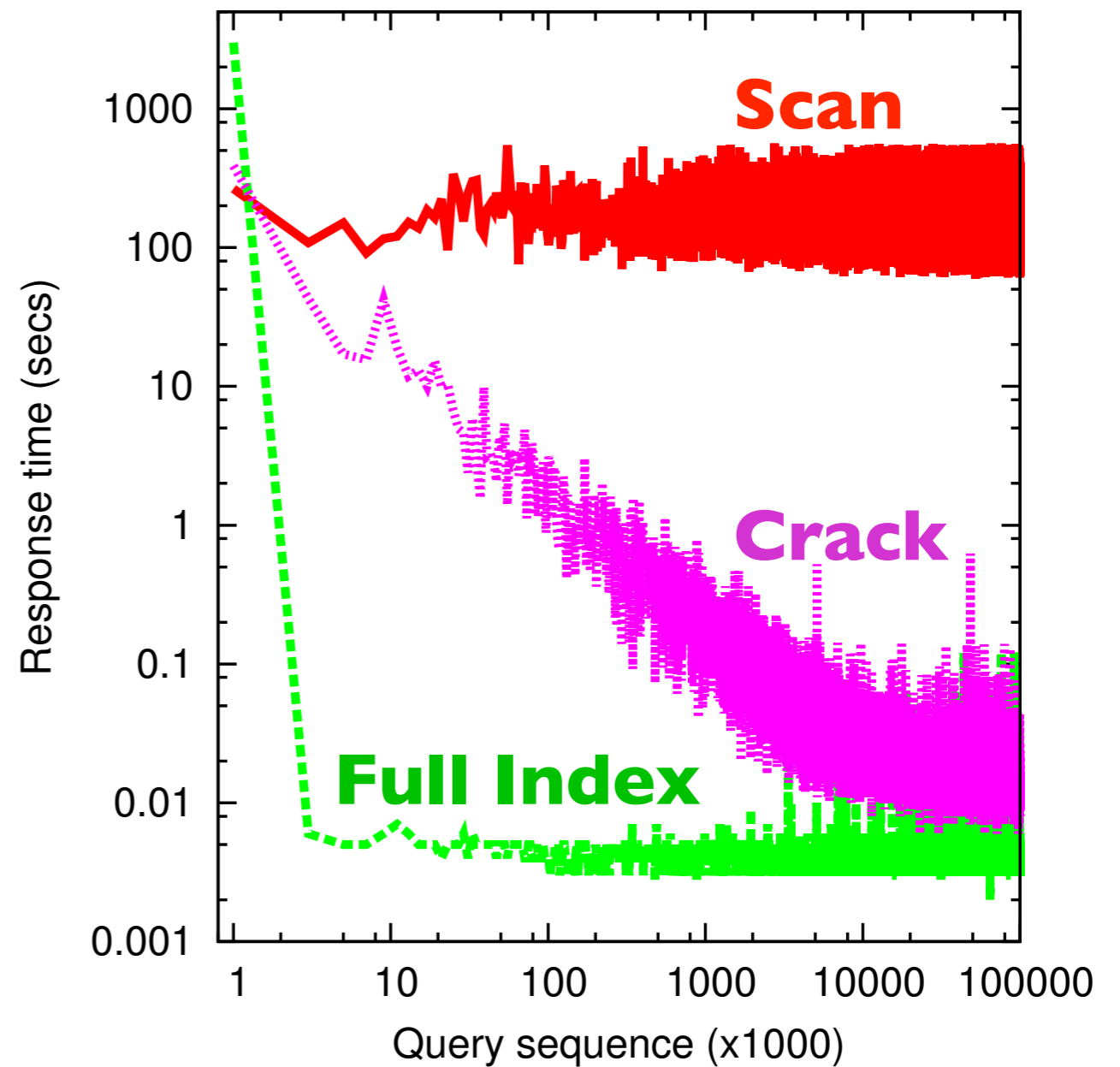


database kernel
code footprint ~2M

continuous adaptation

set-up

100K random selections
random selectivity
random value ranges
in a 10 million integer column

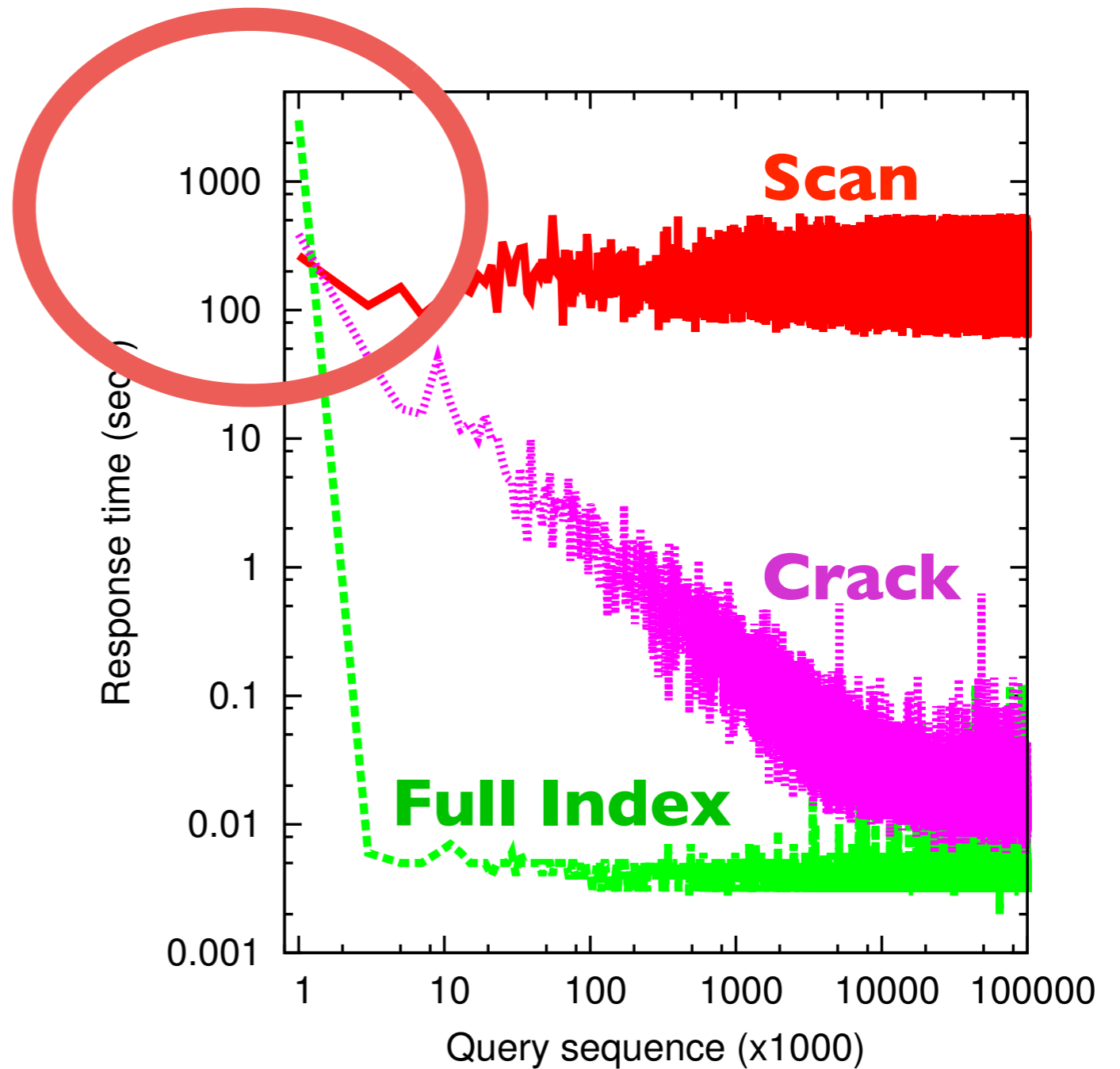


continuous adaptation

set-up

100K random selections
random selectivity
random value ranges
in a 10 million integer column

almost no
initialization overhead

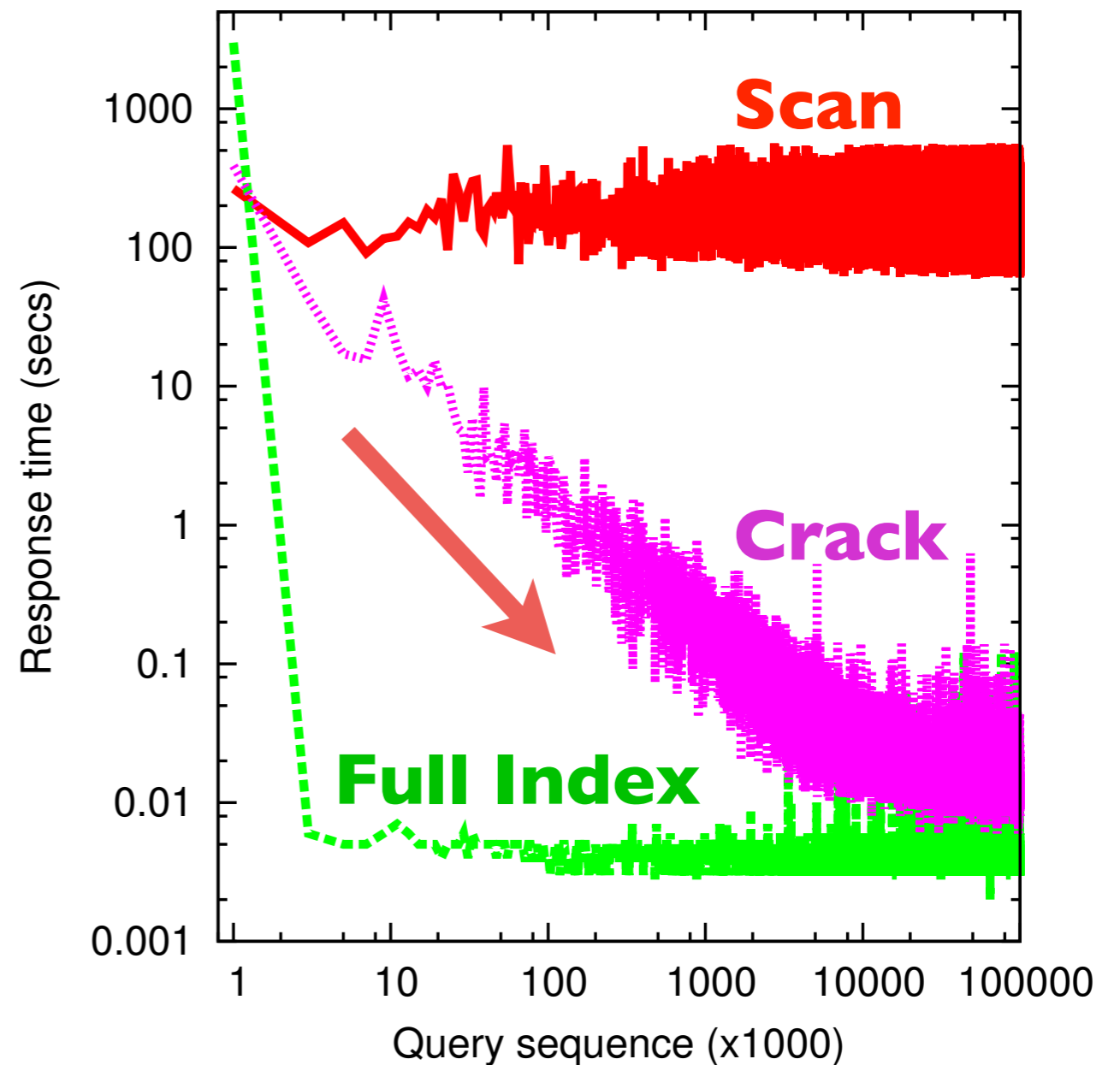


continuous adaptation

set-up

100K random selections
random selectivity
random value ranges
in a 10 million integer column

almost no
initialization overhead
continuous improvement

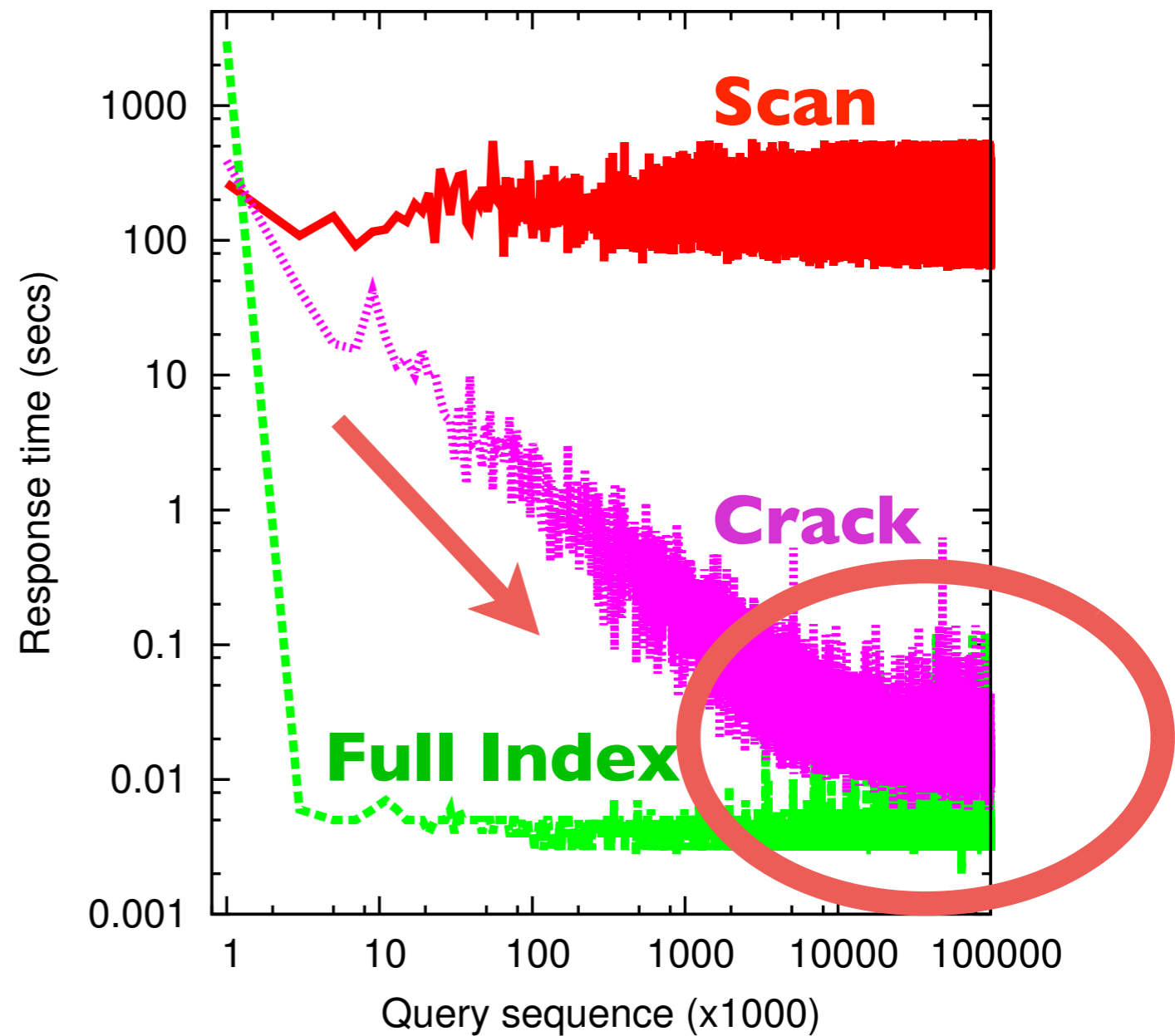


continuous adaptation

set-up

100K random selections
random selectivity
random value ranges
in a 10 million integer column

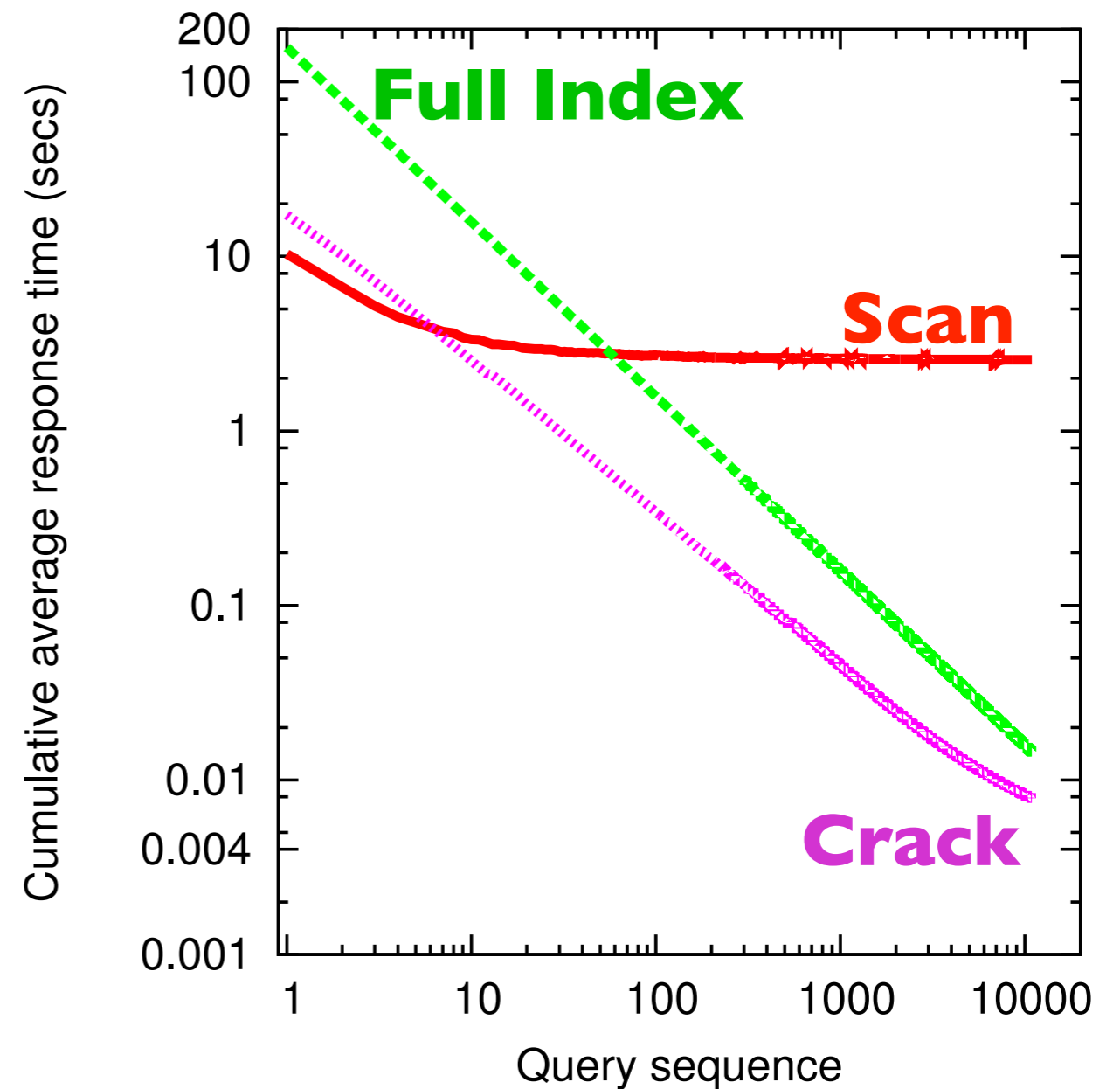
almost no
initialization overhead
continuous improvement



continuous adaptation

set-up

10K random selections
selectivity 10%
random value ranges
in a 30 million integer column

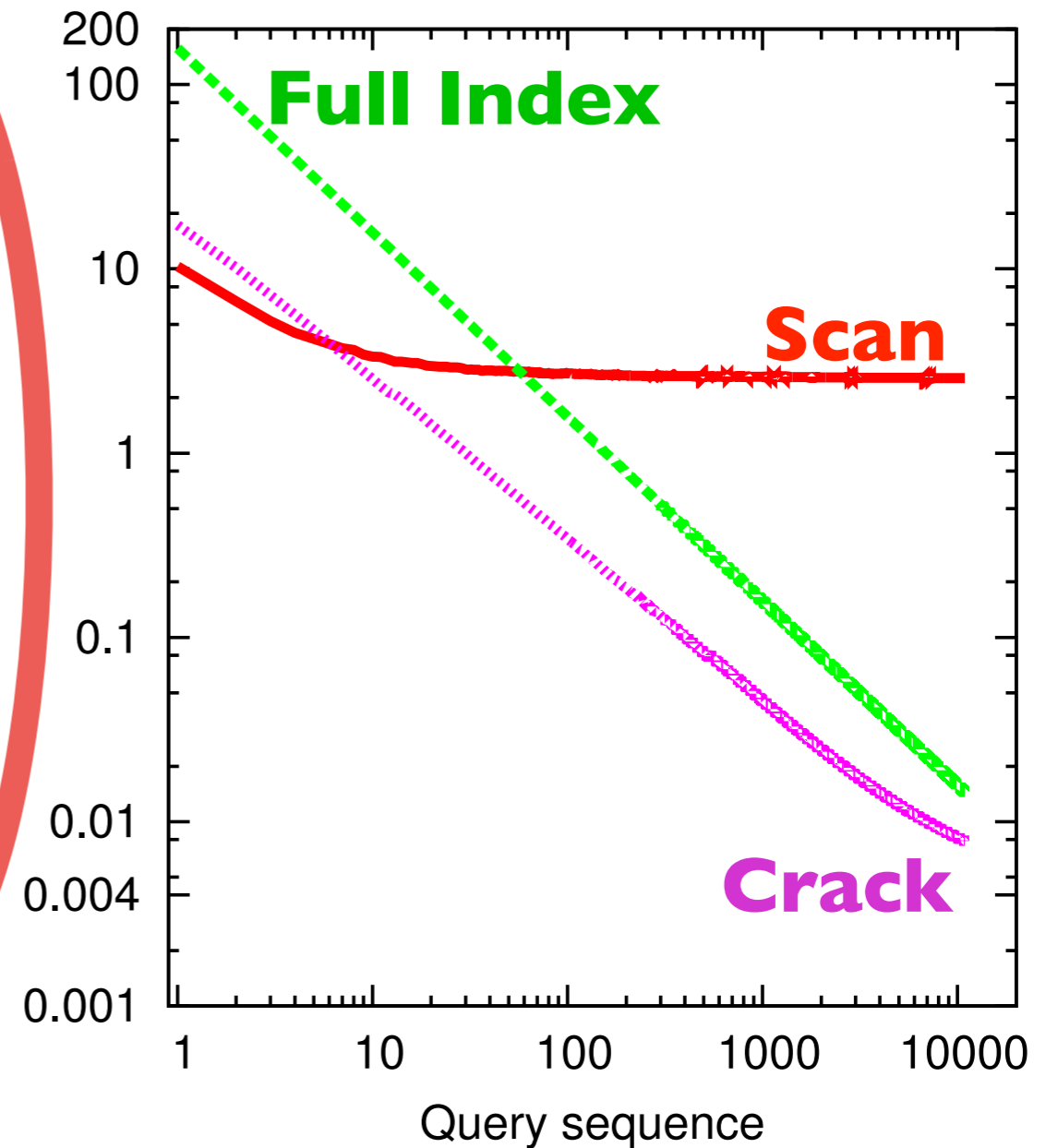


continuous adaptation

set-up

10K random selections
selectivity 10%
random value ranges
in a 30 million integer column

Cumulative average response time (secs)

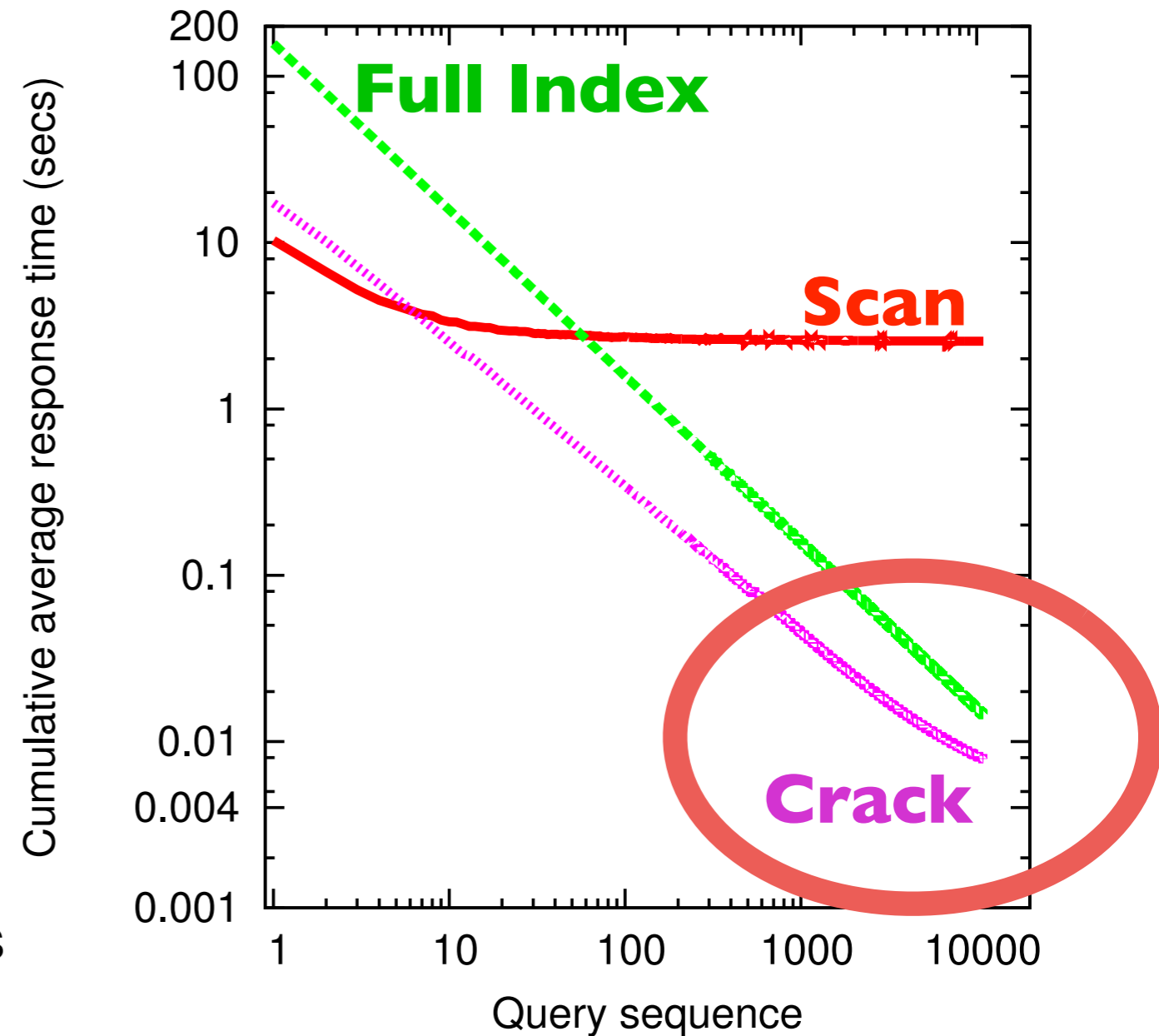


continuous adaptation

set-up

10K random selections
selectivity 10%
random value ranges
in a 30 million integer column

10K queries later,
Full Index still has not
amortized the initialization costs

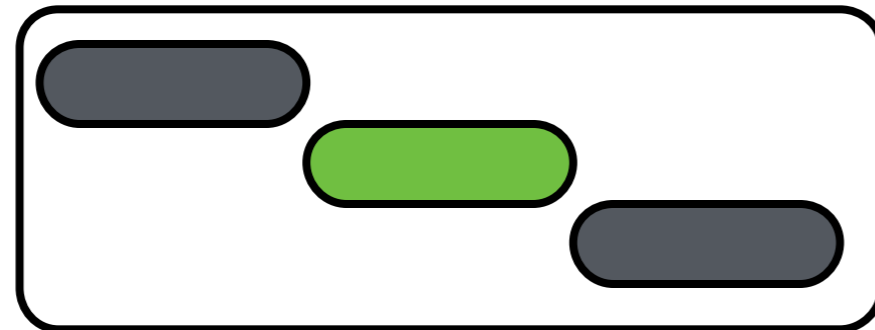


traditional databases

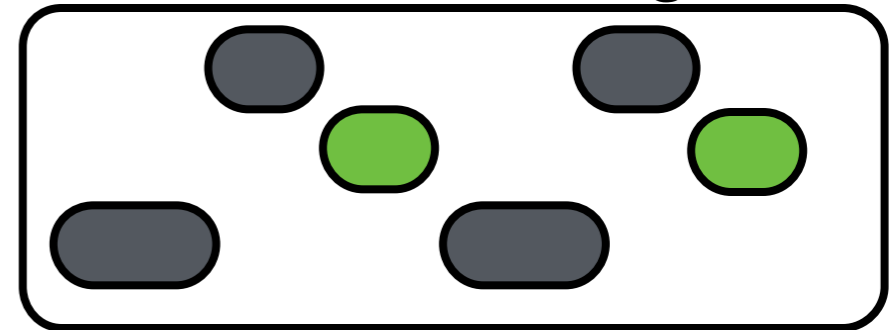
monolithic/full indexing

workload analysis
index building
query processing

offline indexing



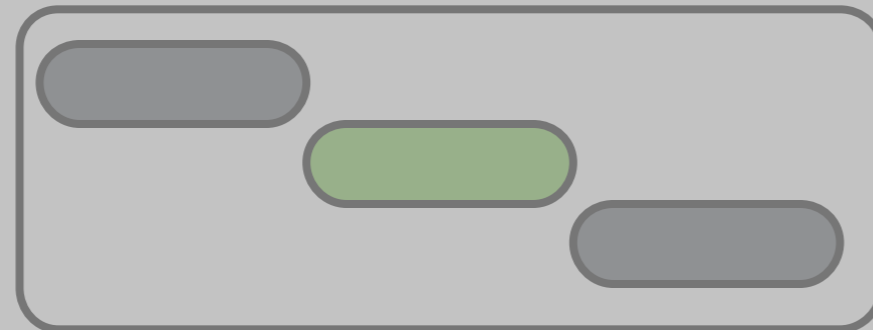
online indexing



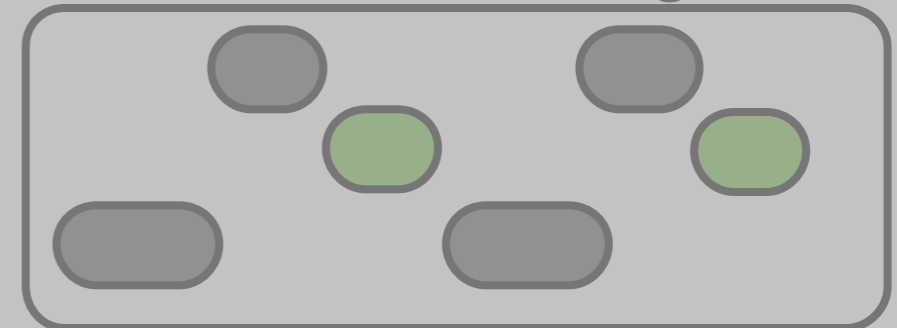
traditional databases

monolithic/full indexing

offline indexing



online indexing



workload analysis
index building
query processing

database cracking

partial/adaptive/continuous indexing

adaptive indexing



table 1

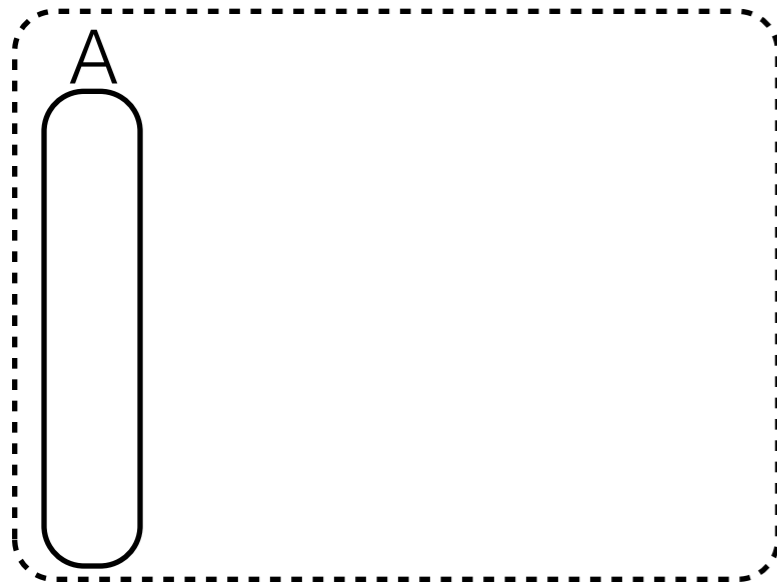
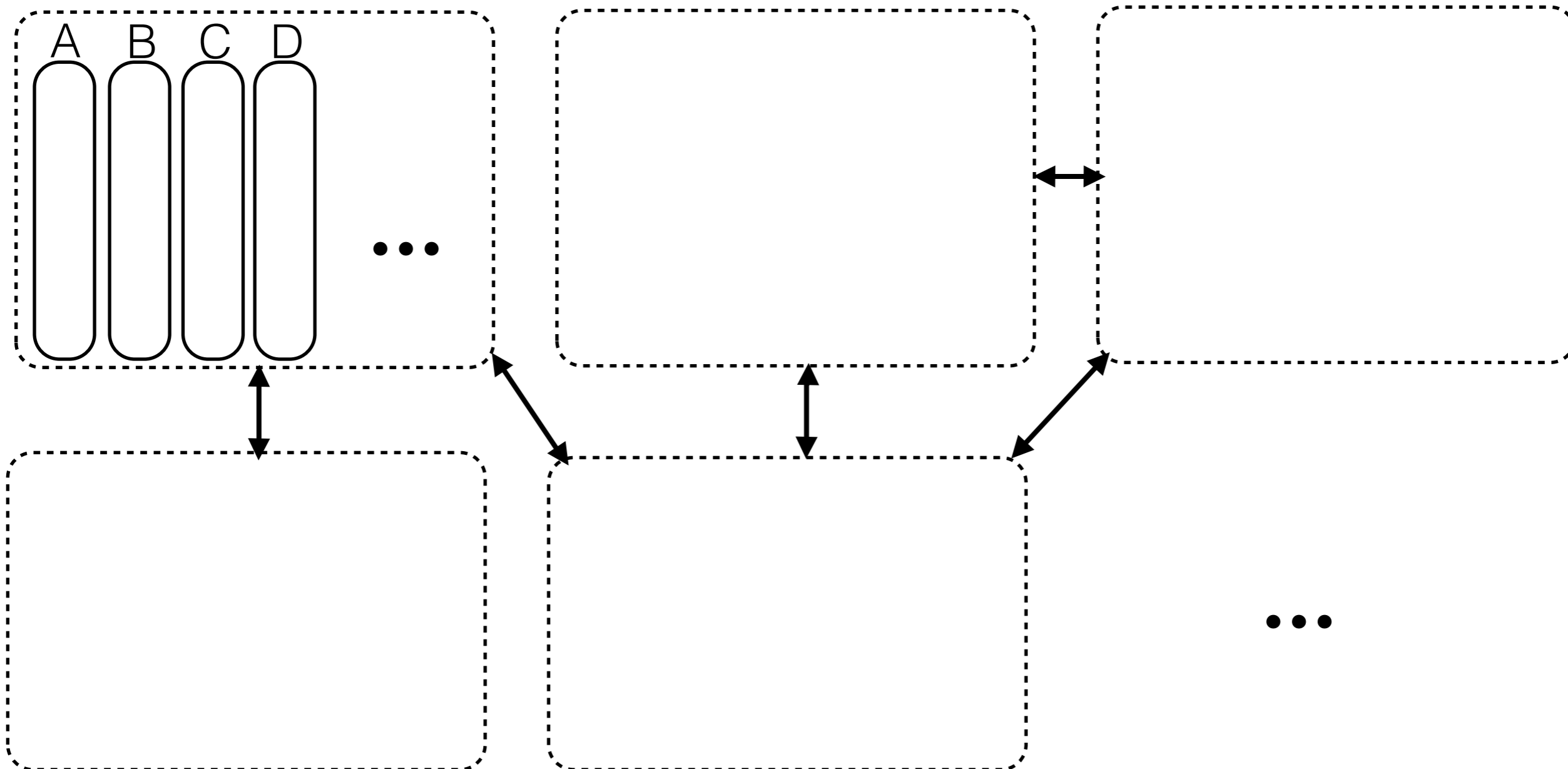
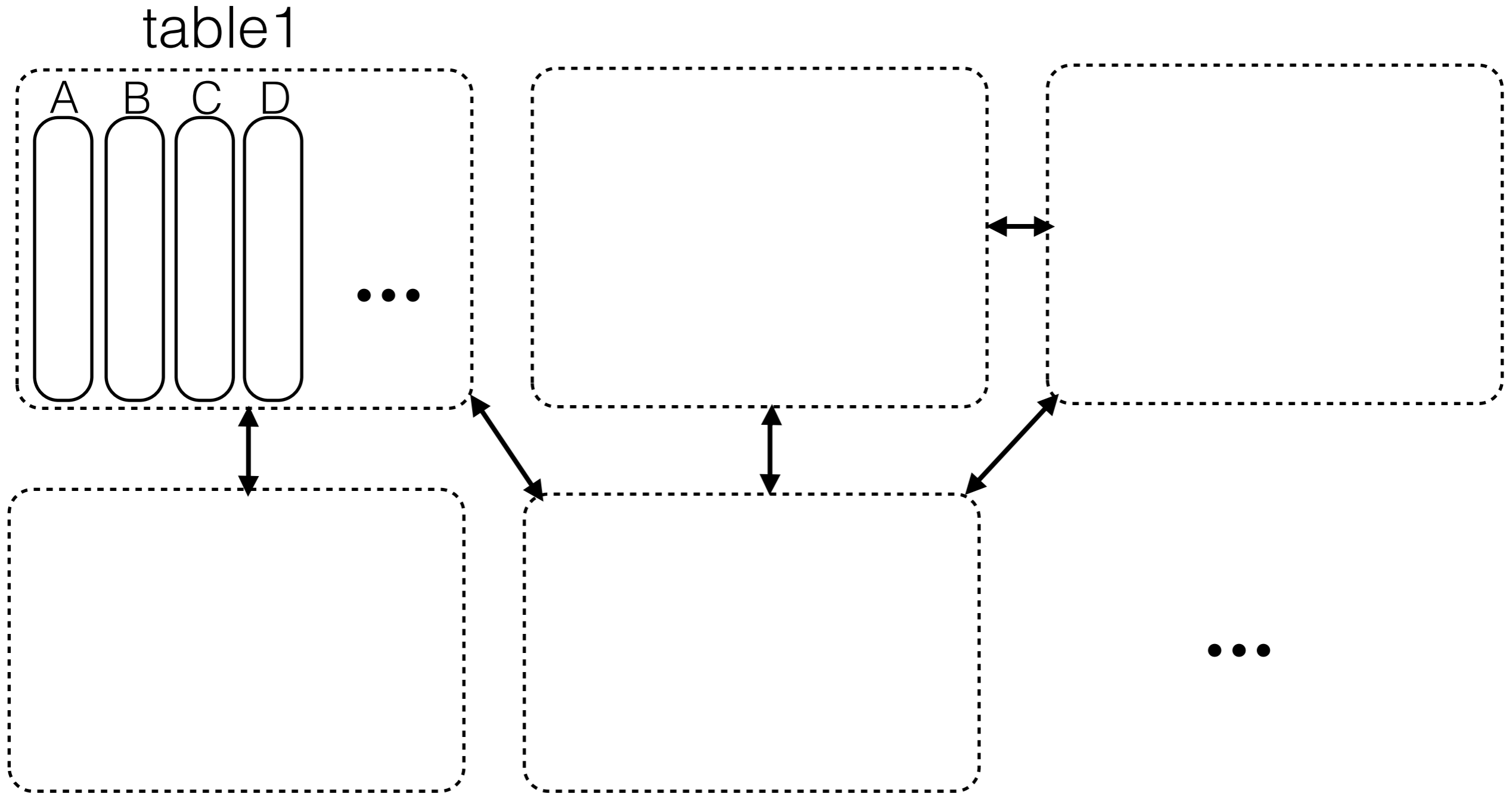


table 1



select R.A from R where R.A > 10 and R.A < 14



select R.A from R where R.A > 10 and R.A < 14

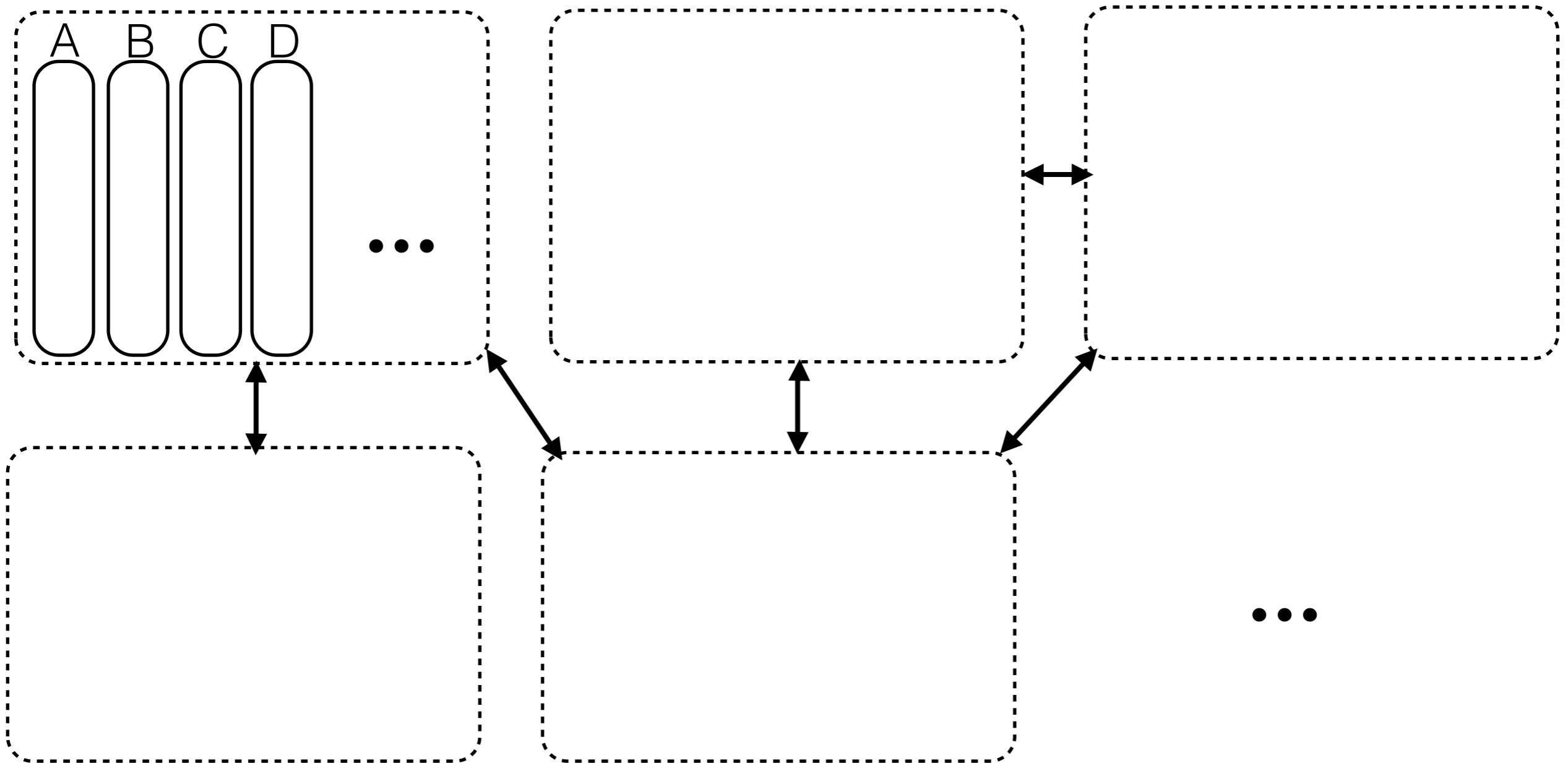
select max(R.A), max(R.B), max(S.A), max(S.B) from R, S

where $v1 < R.C < v2$ and $v3 < R.D < v4$

and $v5 < R.E < v6$ and $k1 < S.C < k2$ and $k3 < S.D < k4$ and $k5 < S.E < k6$

and R.F = S.F

table 1



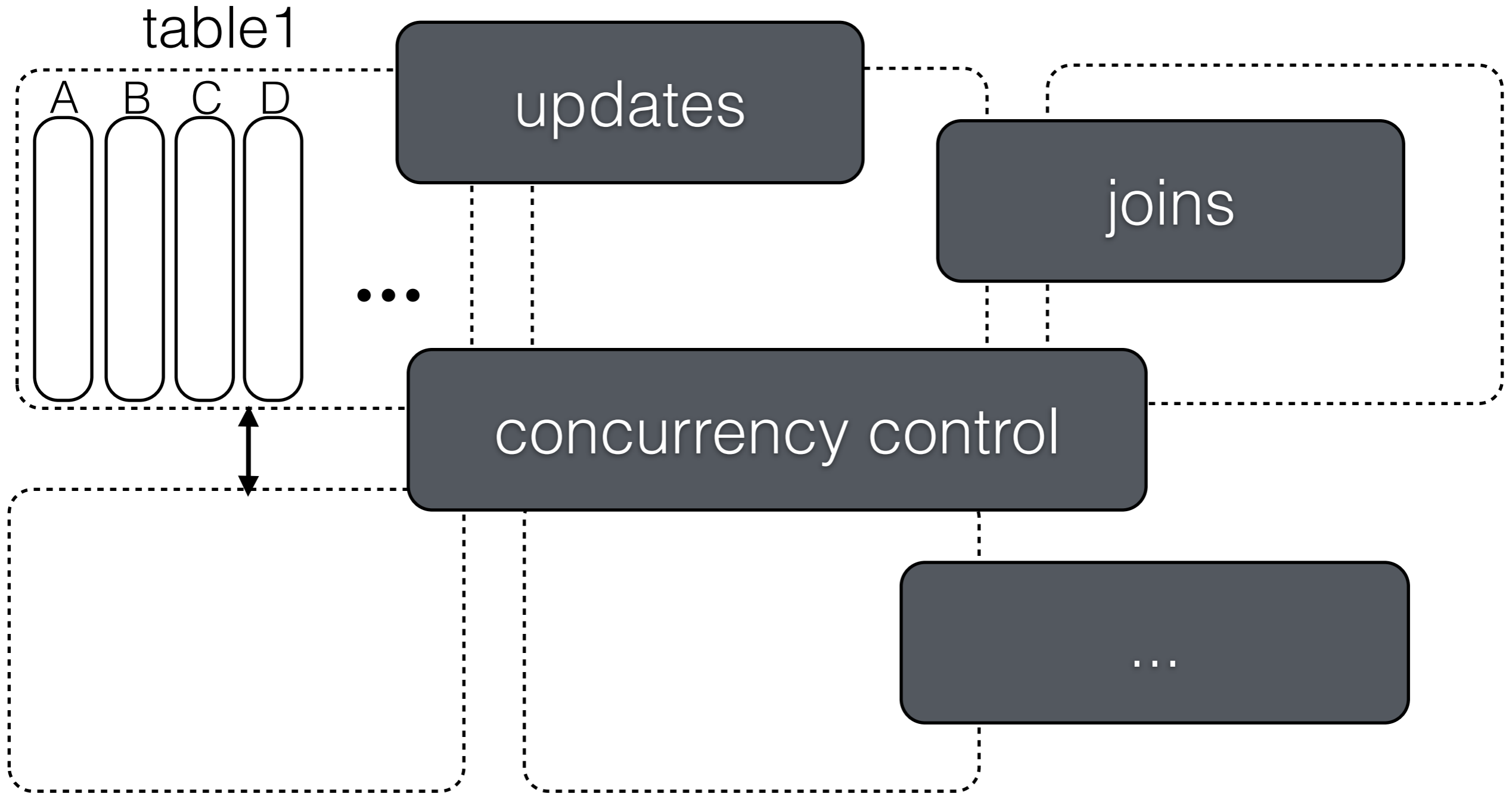
select R.A from R where R.A > 10 and R.A < 14

select max(R.A), max(R.B), max(S.A), max(S.B) from R, S

where v1 < R.C < v2 and v3 < R.D < v4

and v5 < R.E < v6 and k1 < S.C < k2 and k3 < S.D < k4 and k5 < S.E < k6

and R.F = S.F



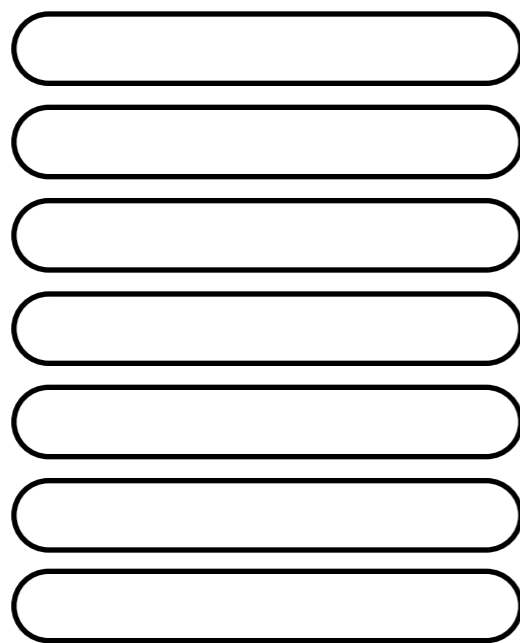
sideways cracking

tuple reconstruction

rows & columns

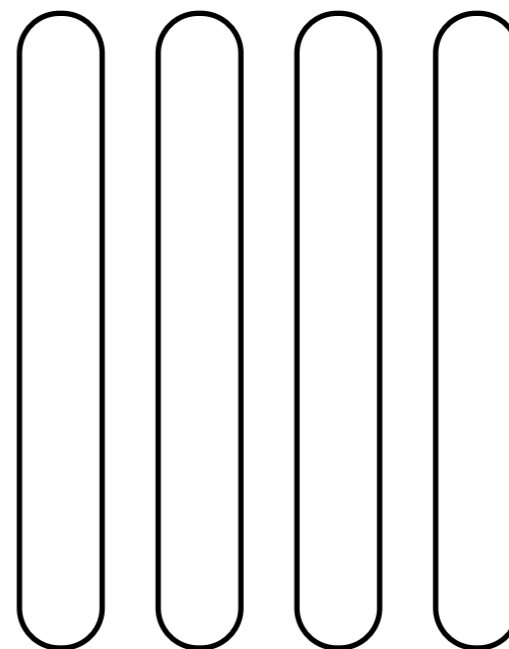
row-store

A B C D

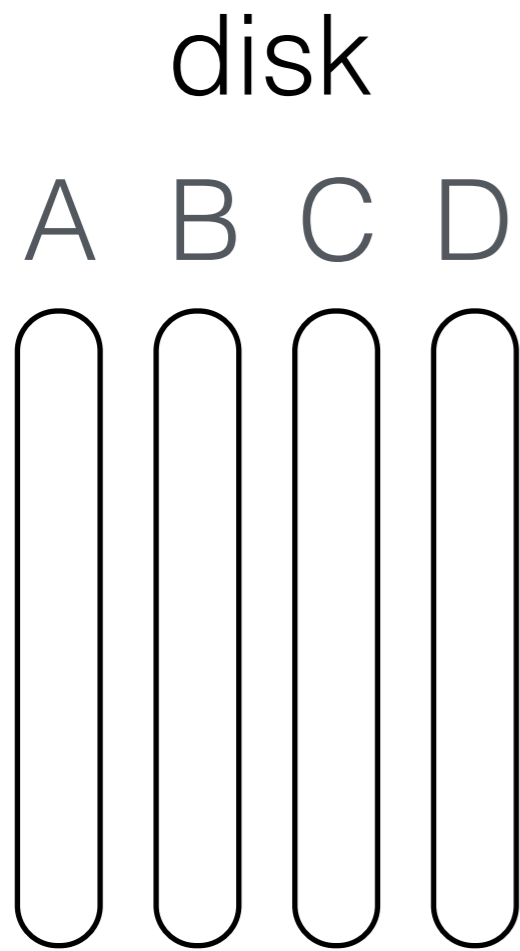


column-store

A B C D

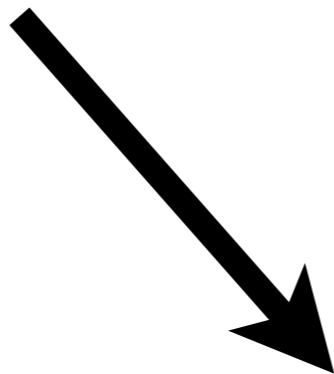
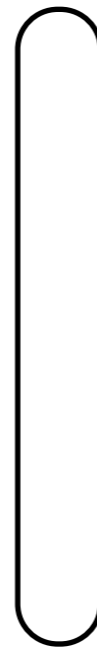


tuple reconstruction



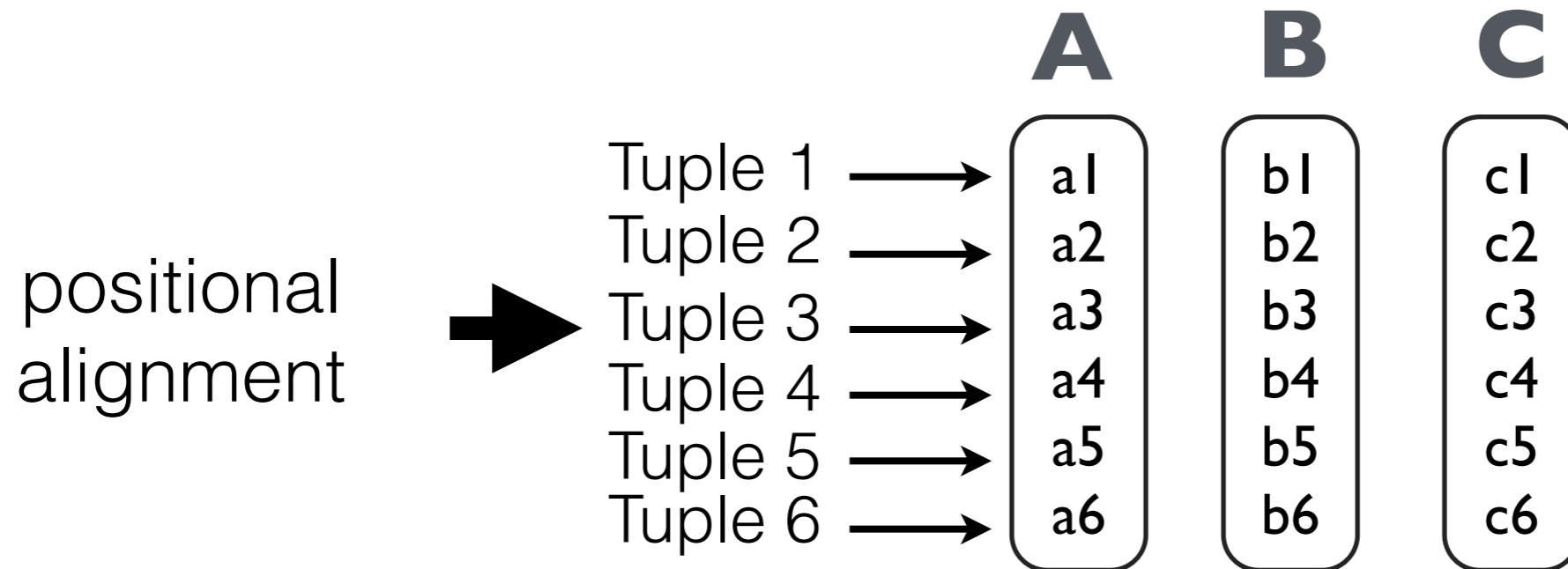
memory

A



A B C D





positional lookups

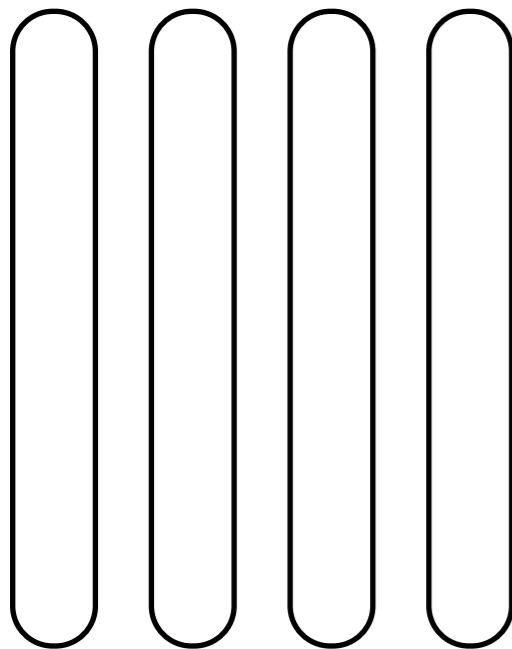
$$A(i) = A + i * \text{width}(A)$$

select min(C) from R where $A < 10$ & $B < 20$



disk

A **B** **C** **D**

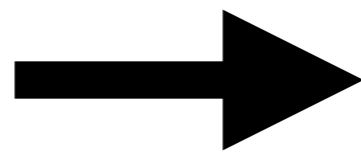
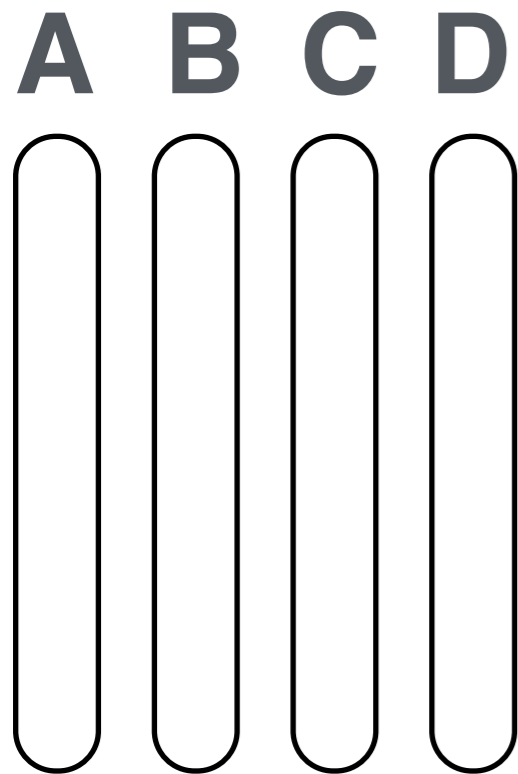


select min(C) from R where $A < 10$ & $B < 20$

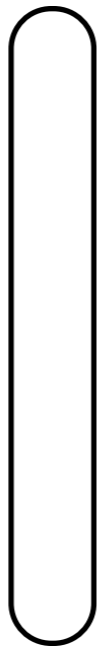


disk

memory



$A < 10$

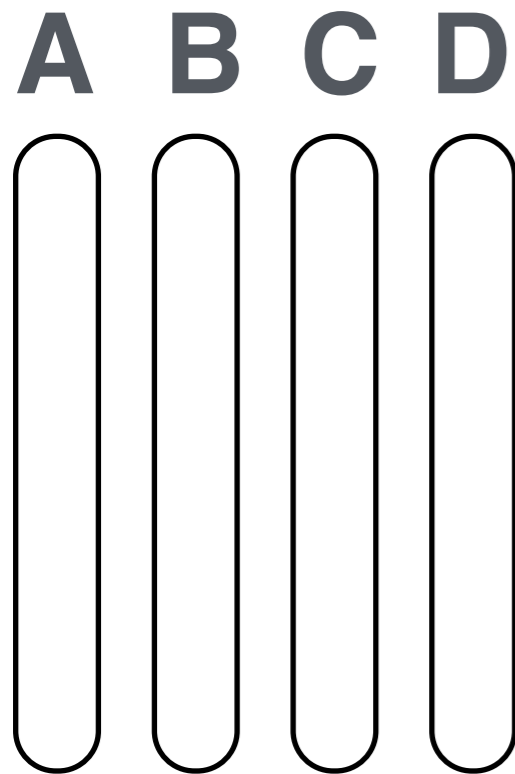


select min(C) from R where A < 10 & B < 20

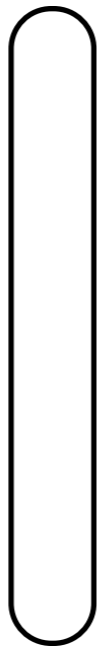
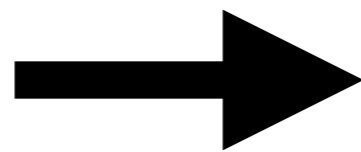


disk

memory

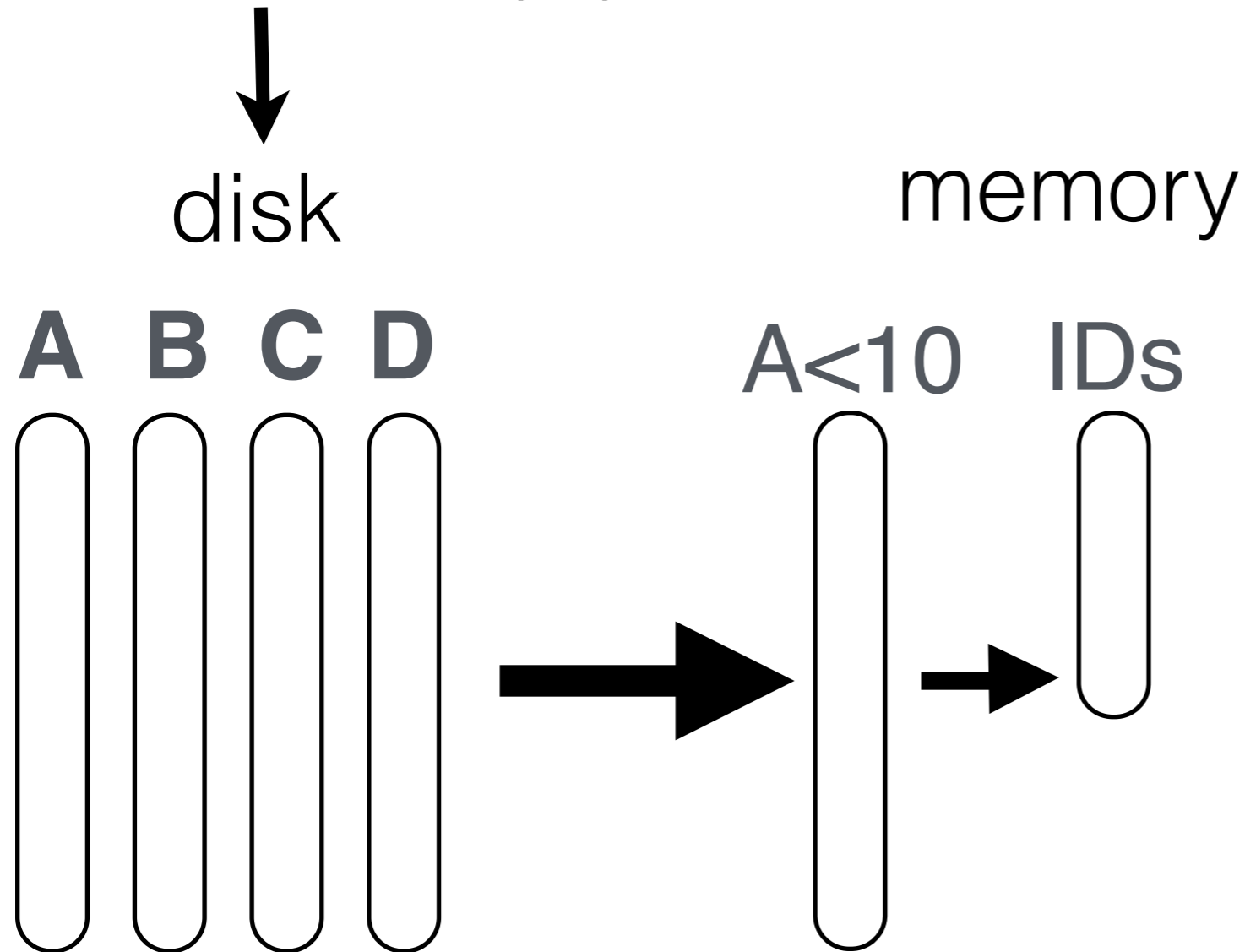


A < 10

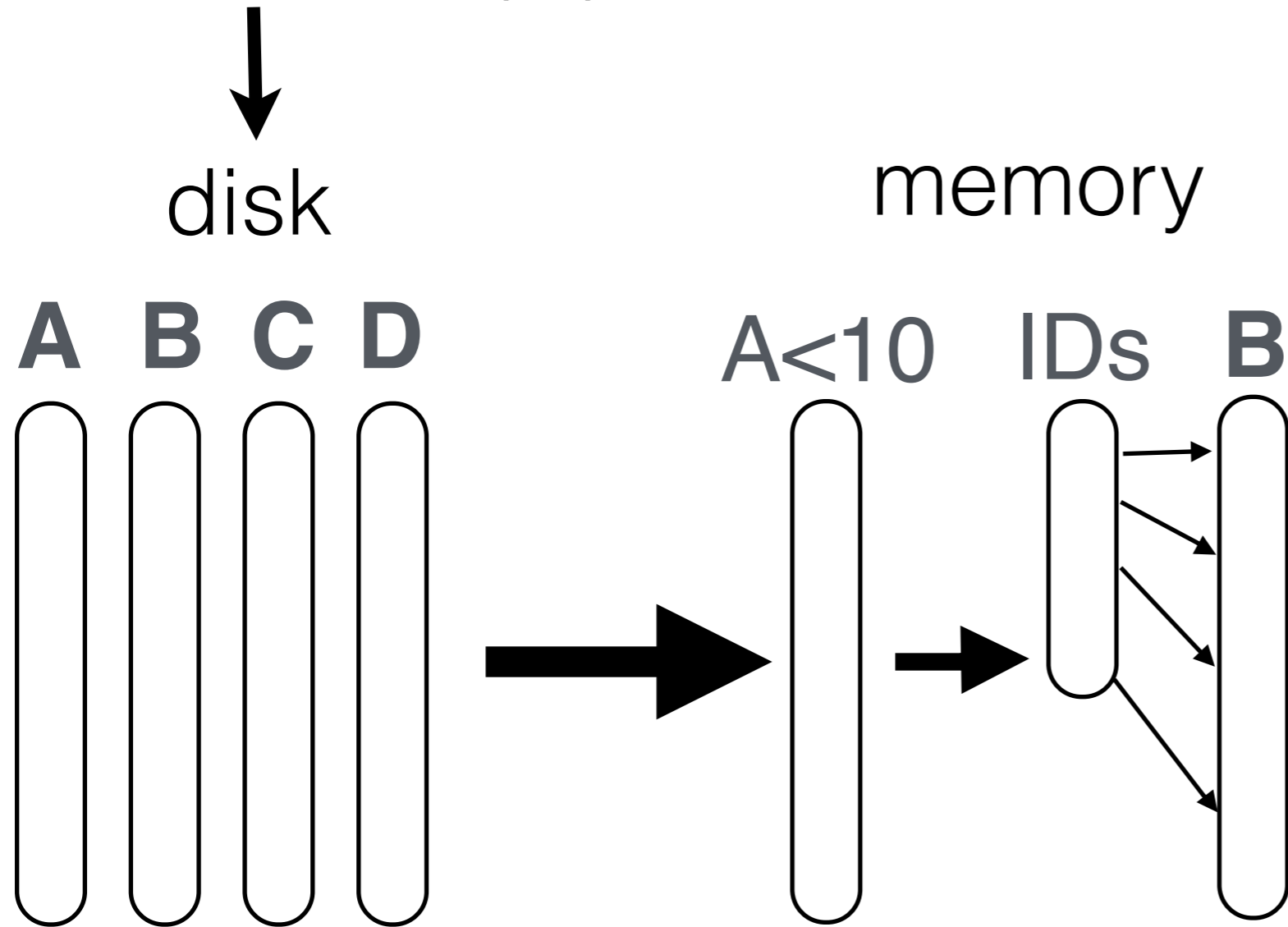


```
1: int *input=A          A < 10
2: for (i=0;i<tuples;i++,input++)
3:     if *input<10
4:         *output=i
5:         output++
```

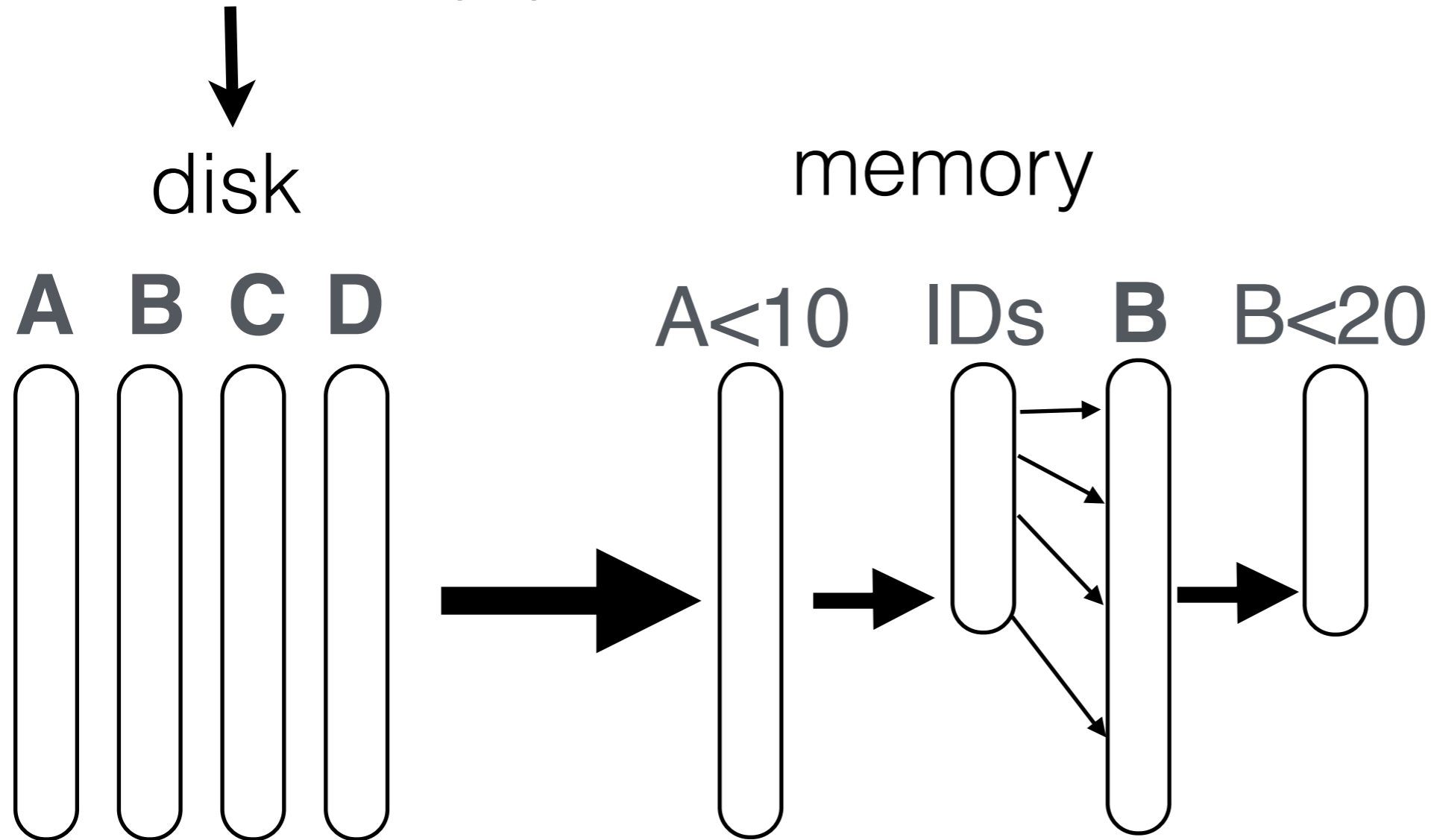
select min(C) from R where $A < 10$ & $B < 20$



select min(C) from R where $A < 10$ & $B < 20$



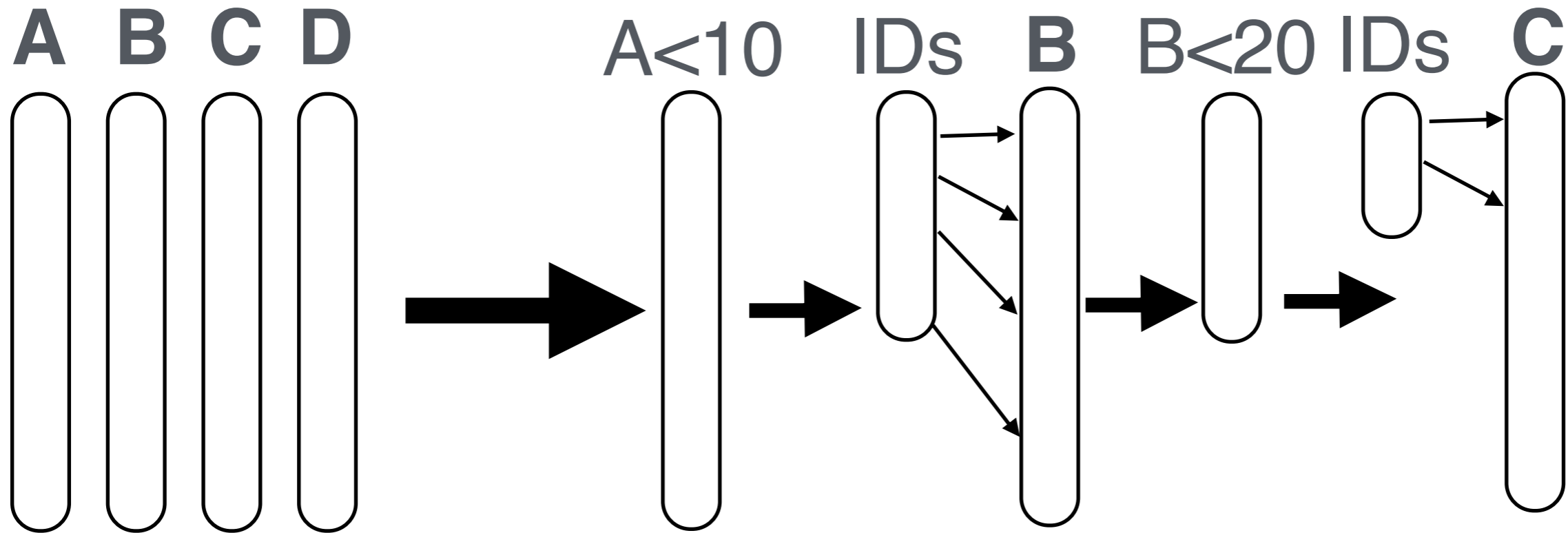
select min(C) from R where $A < 10$ & $B < 20$



select min(C) from R where $A < 10$ & $B < 20$

↓
disk

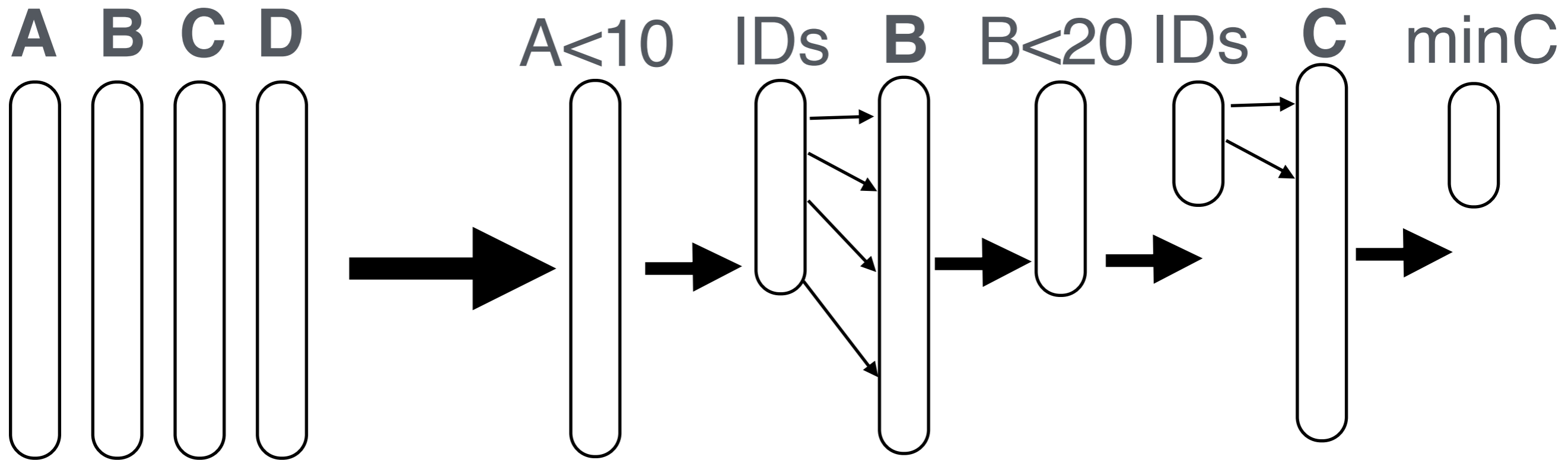
memory



select $\min(C)$ from R where $A < 10$ & $B < 20$

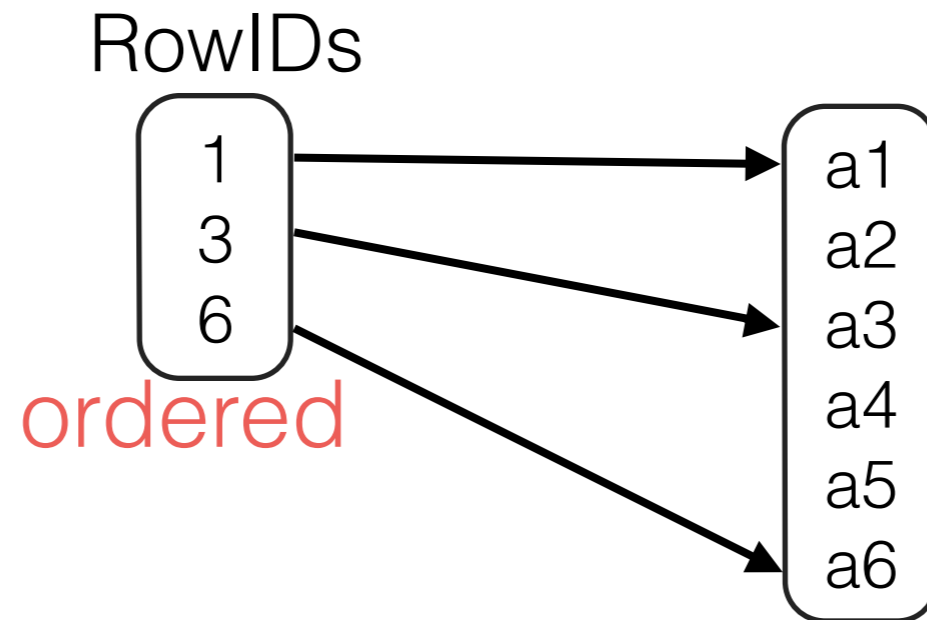
↓
disk

memory



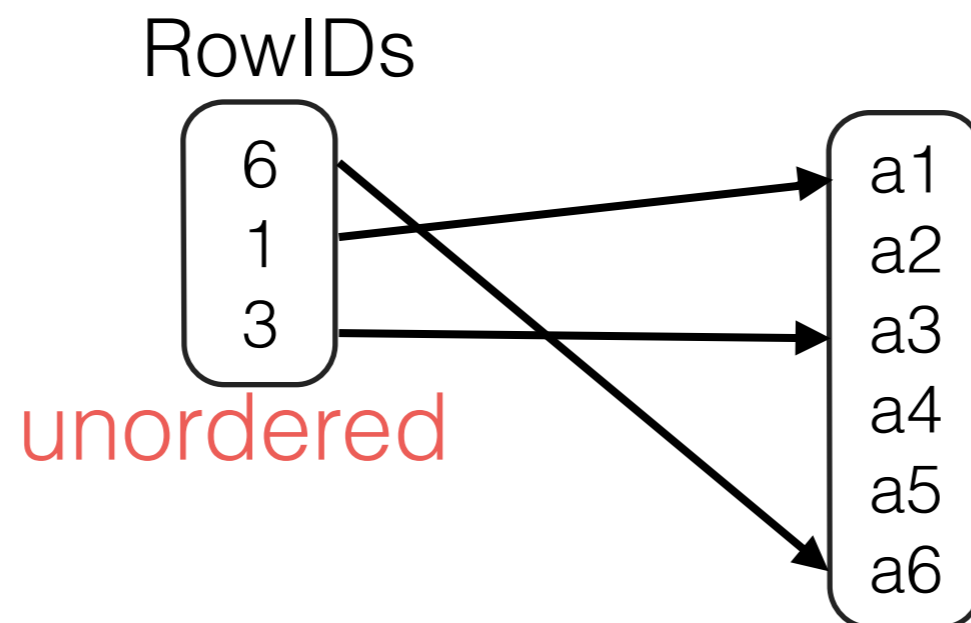
the tuple reconstruction problem

without cracking



sequential access

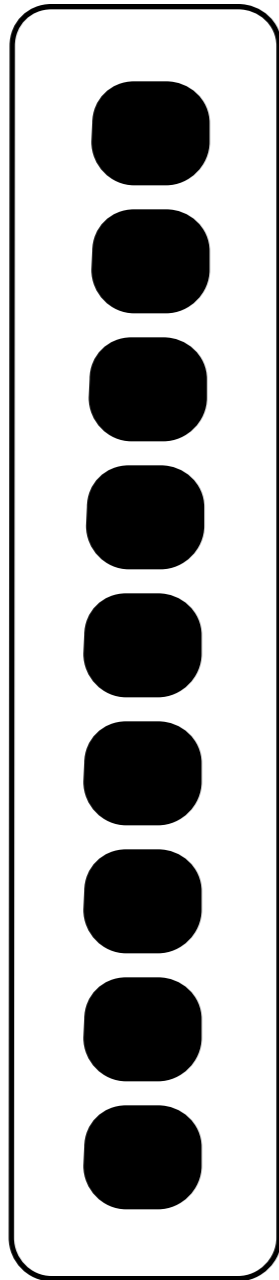
with cracking



random access

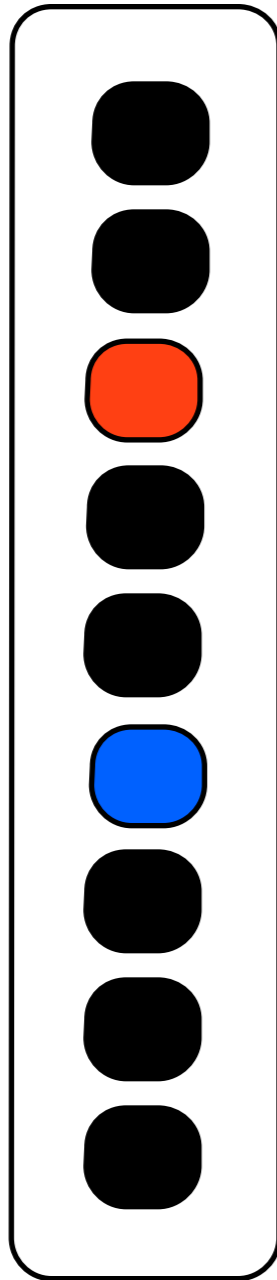
sideways cracking

A

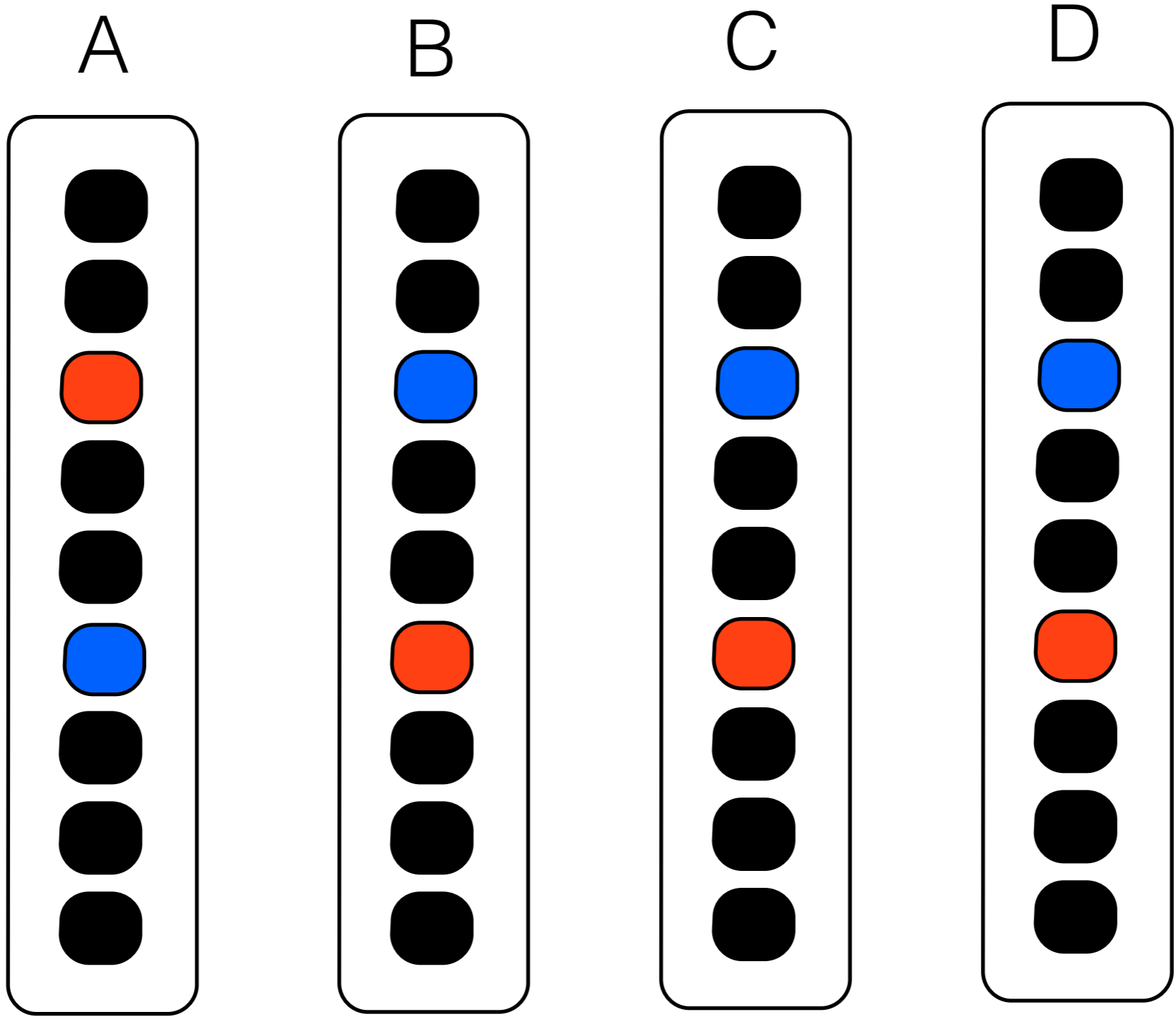


sideways cracking

A

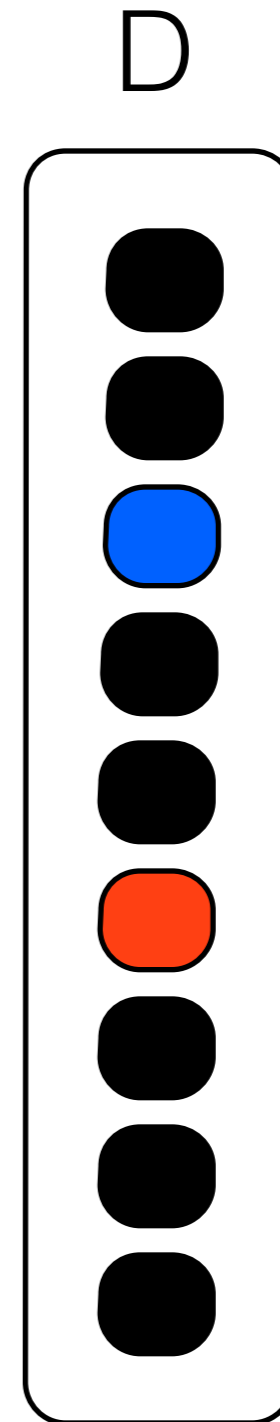
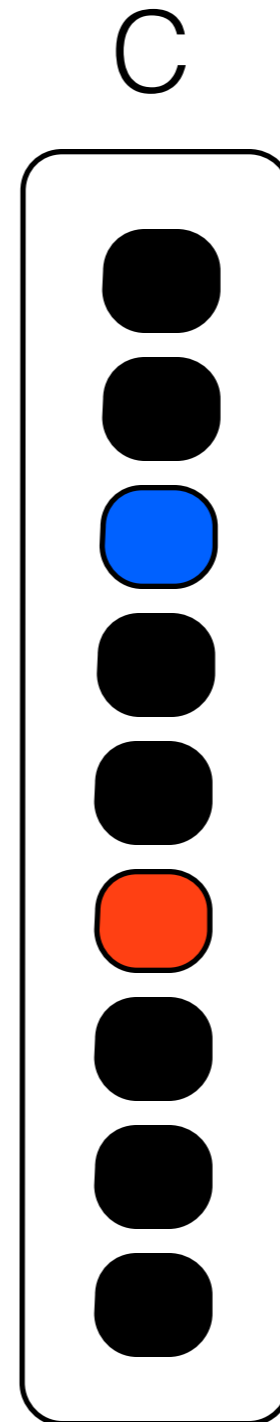
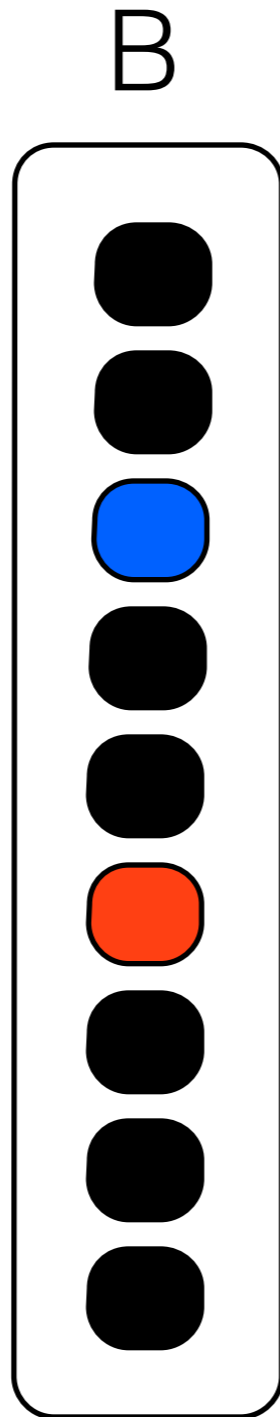
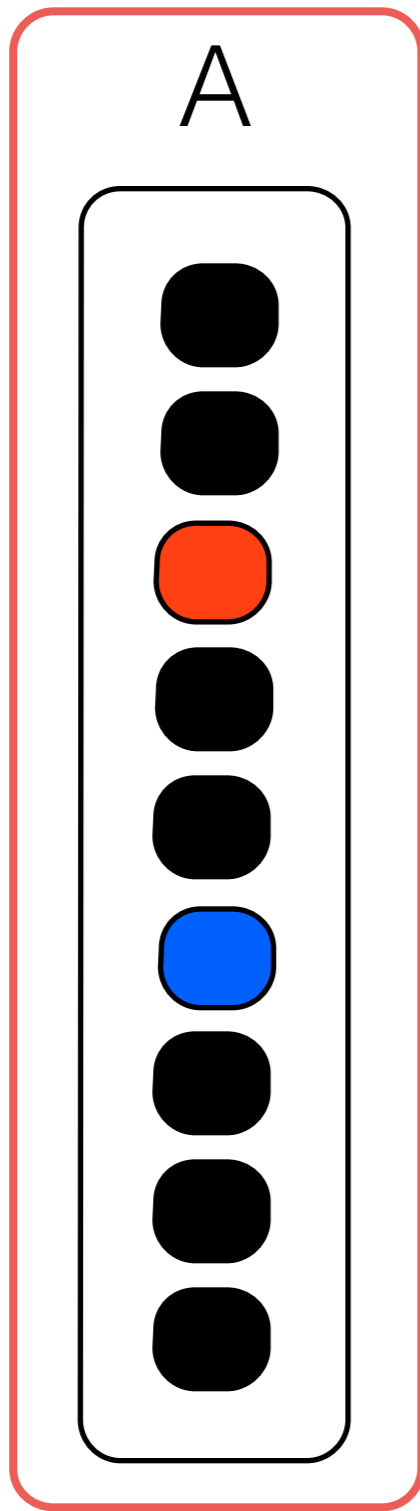


sideways cracking



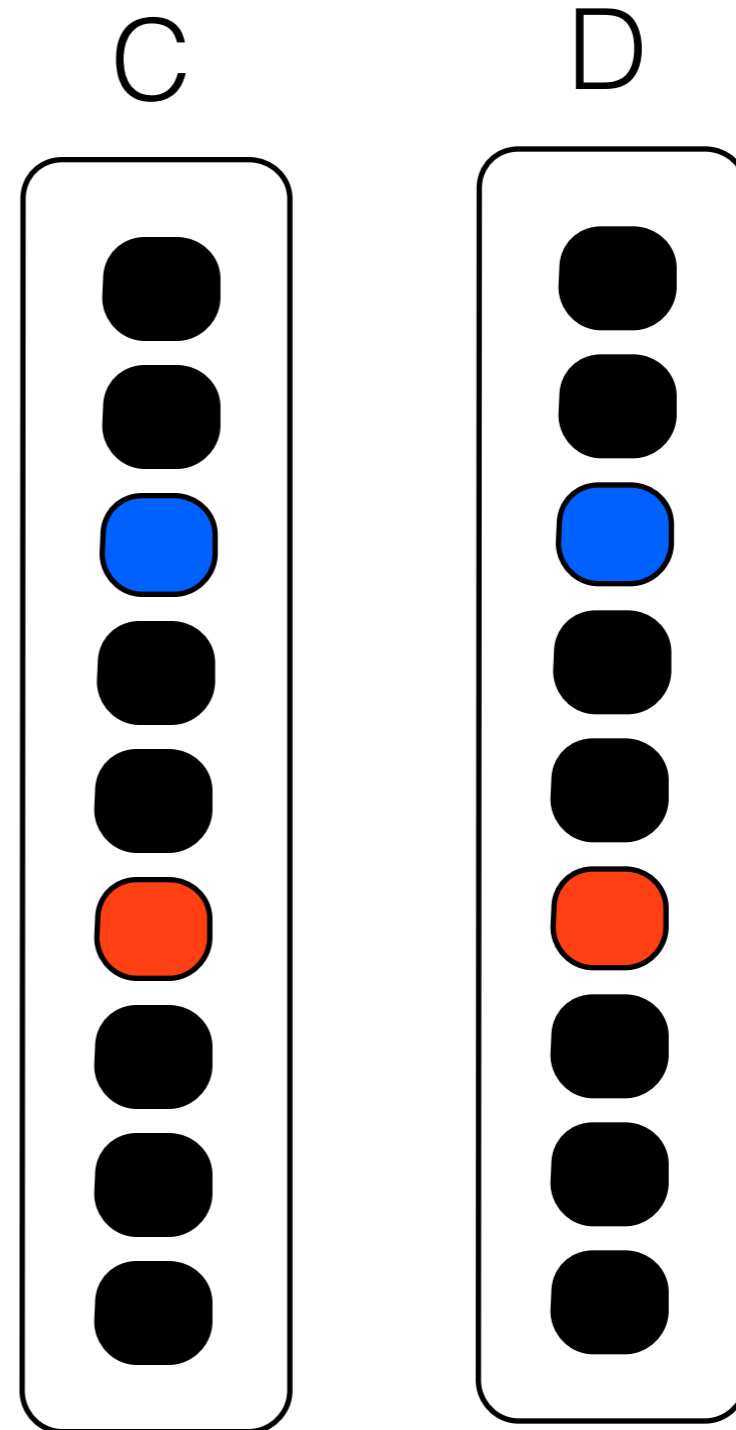
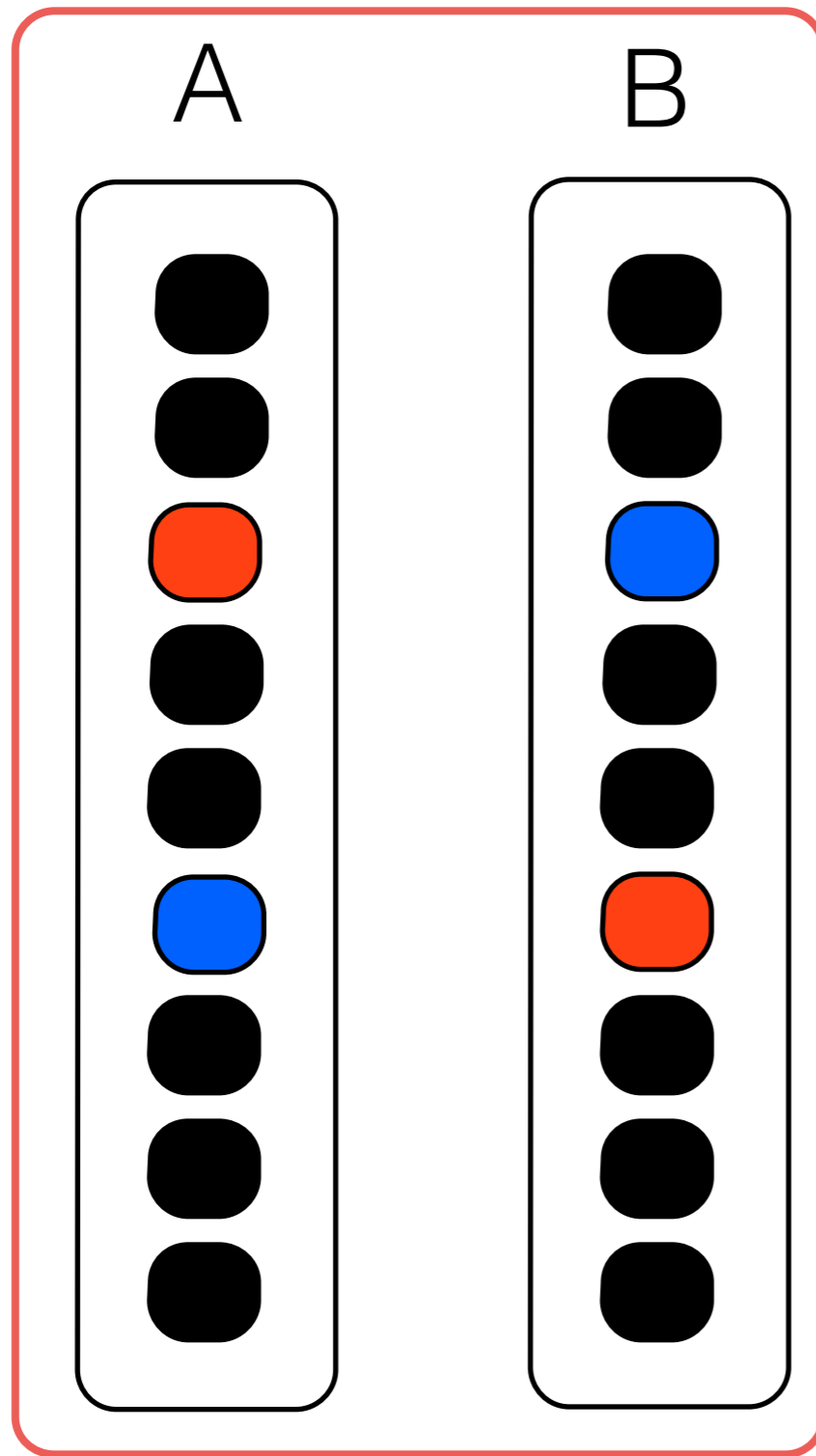
sideways cracking

query



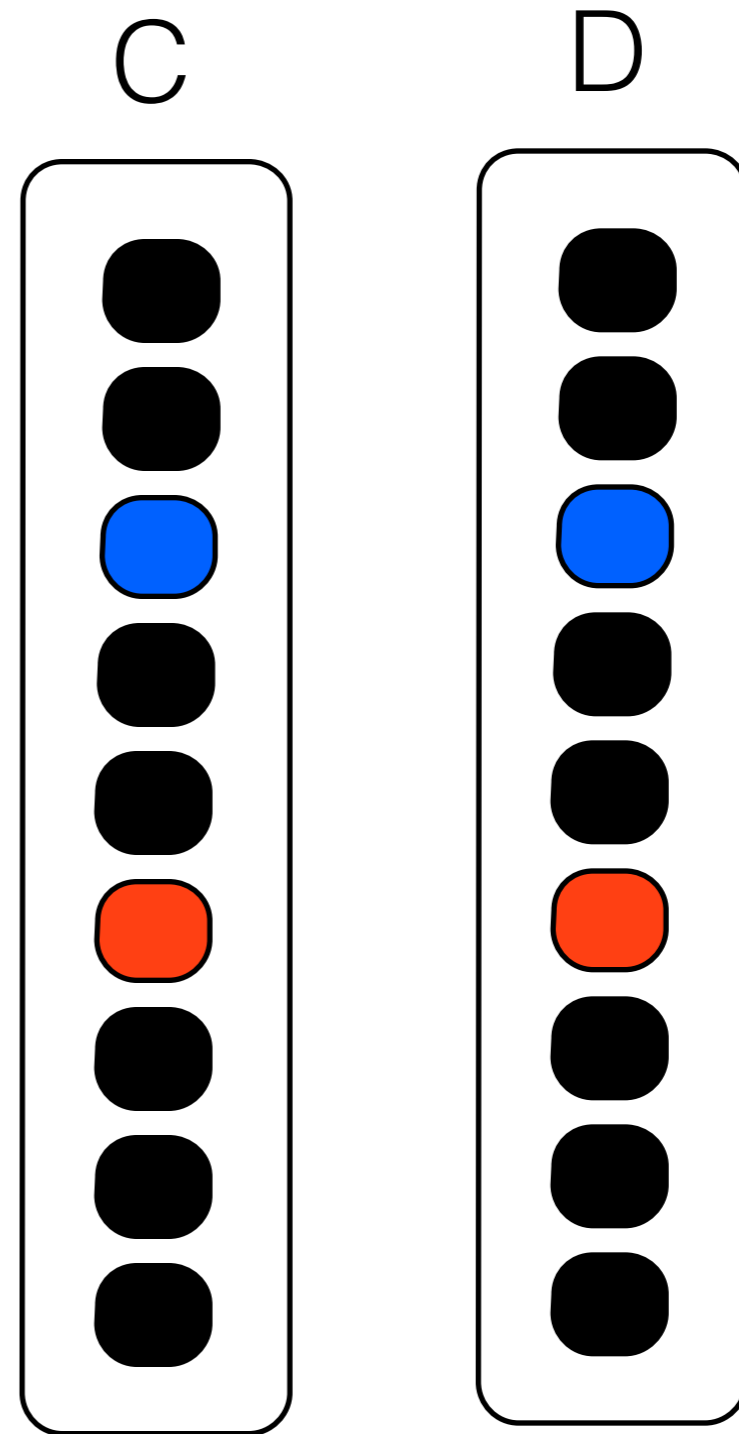
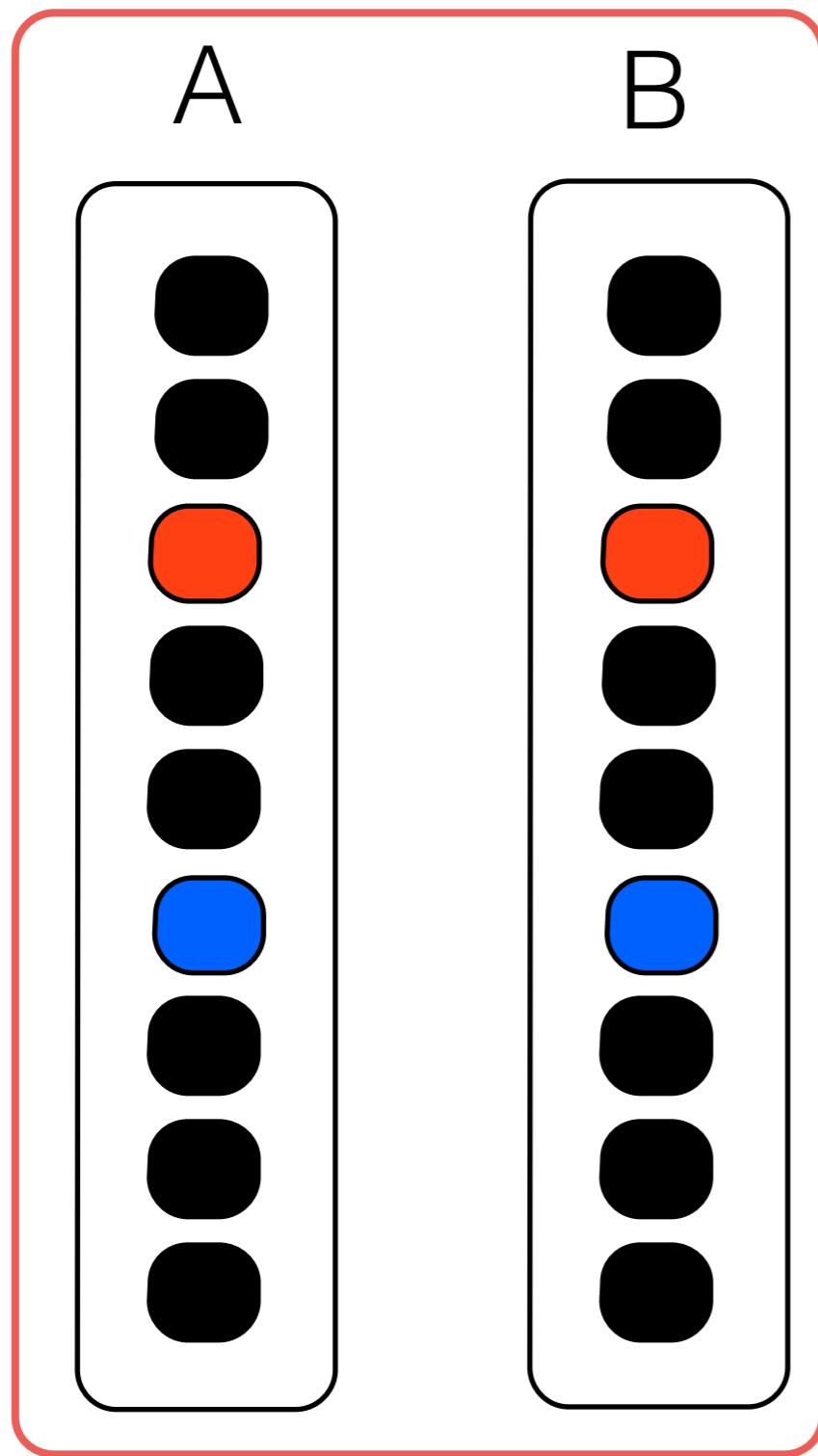
sideways cracking

query



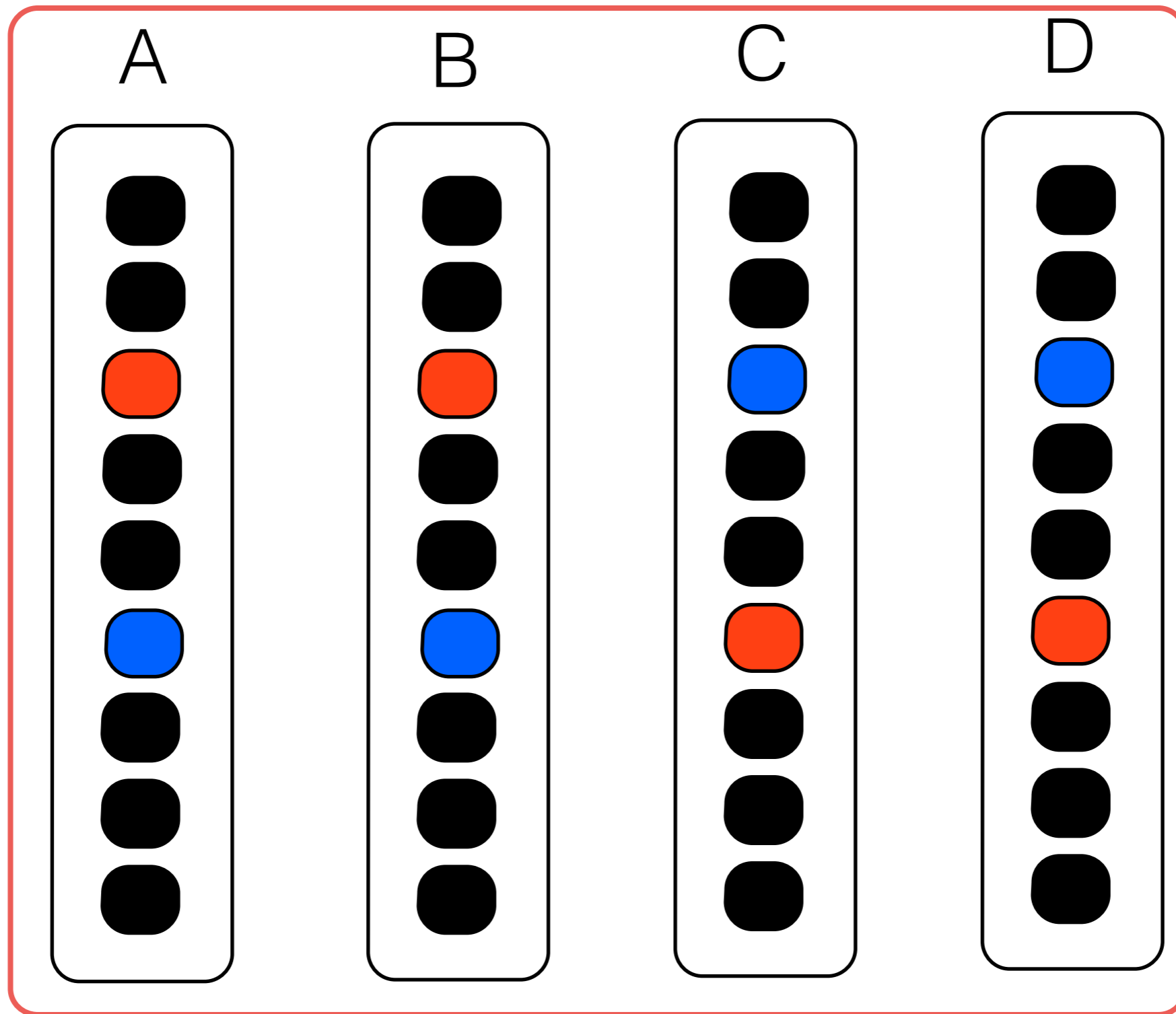
sideways cracking

query



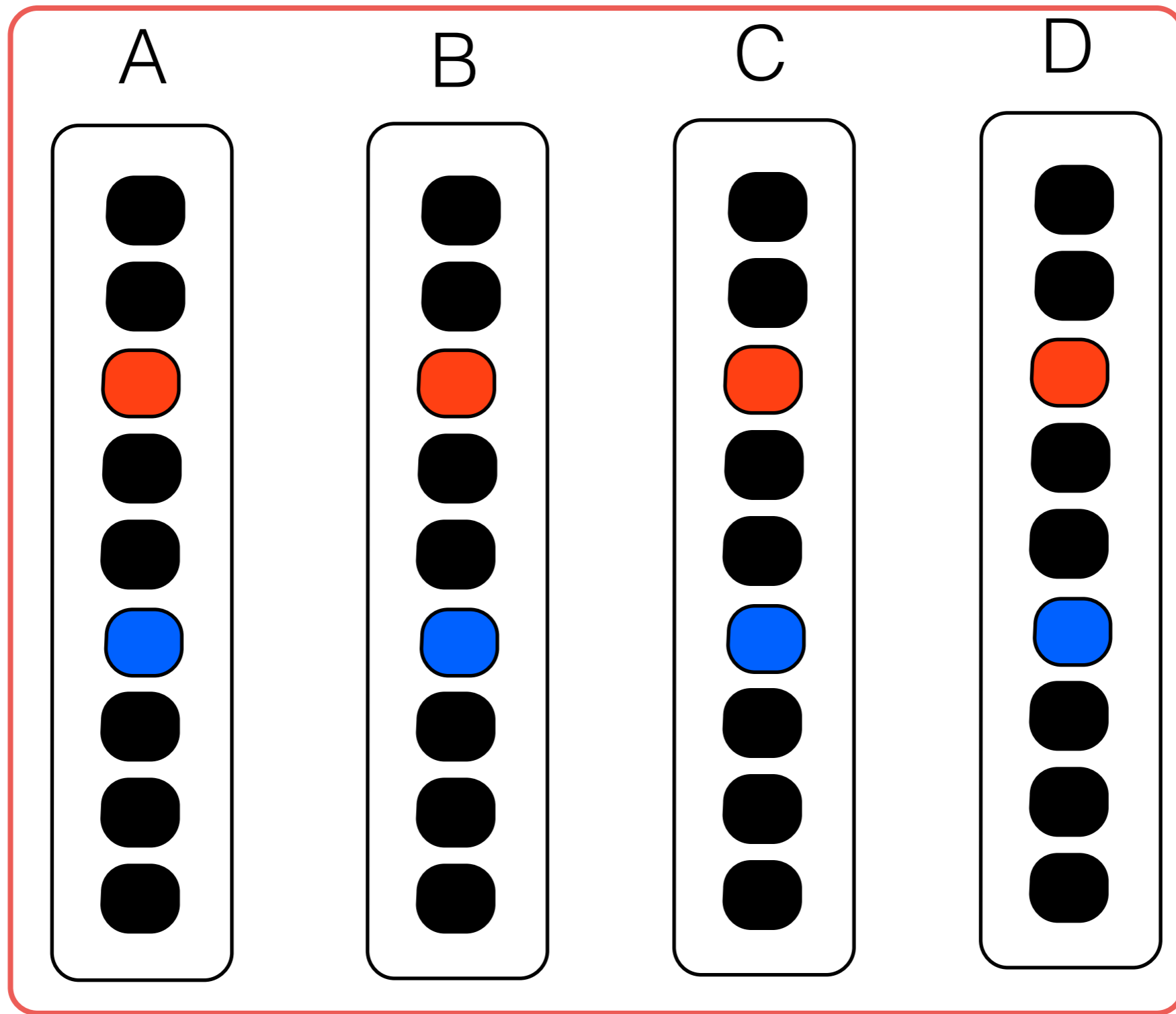
sideways cracking

query

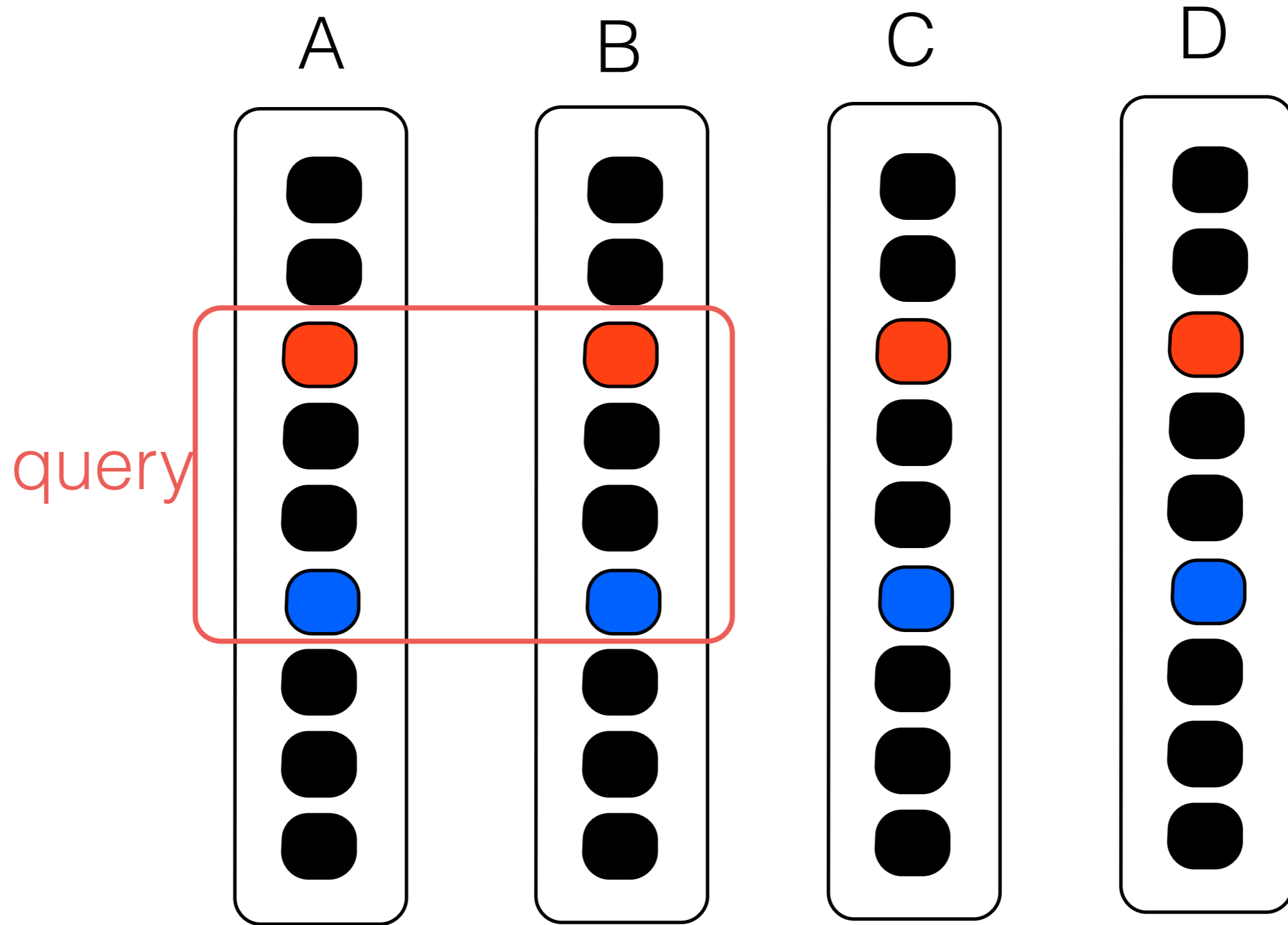


sideways cracking

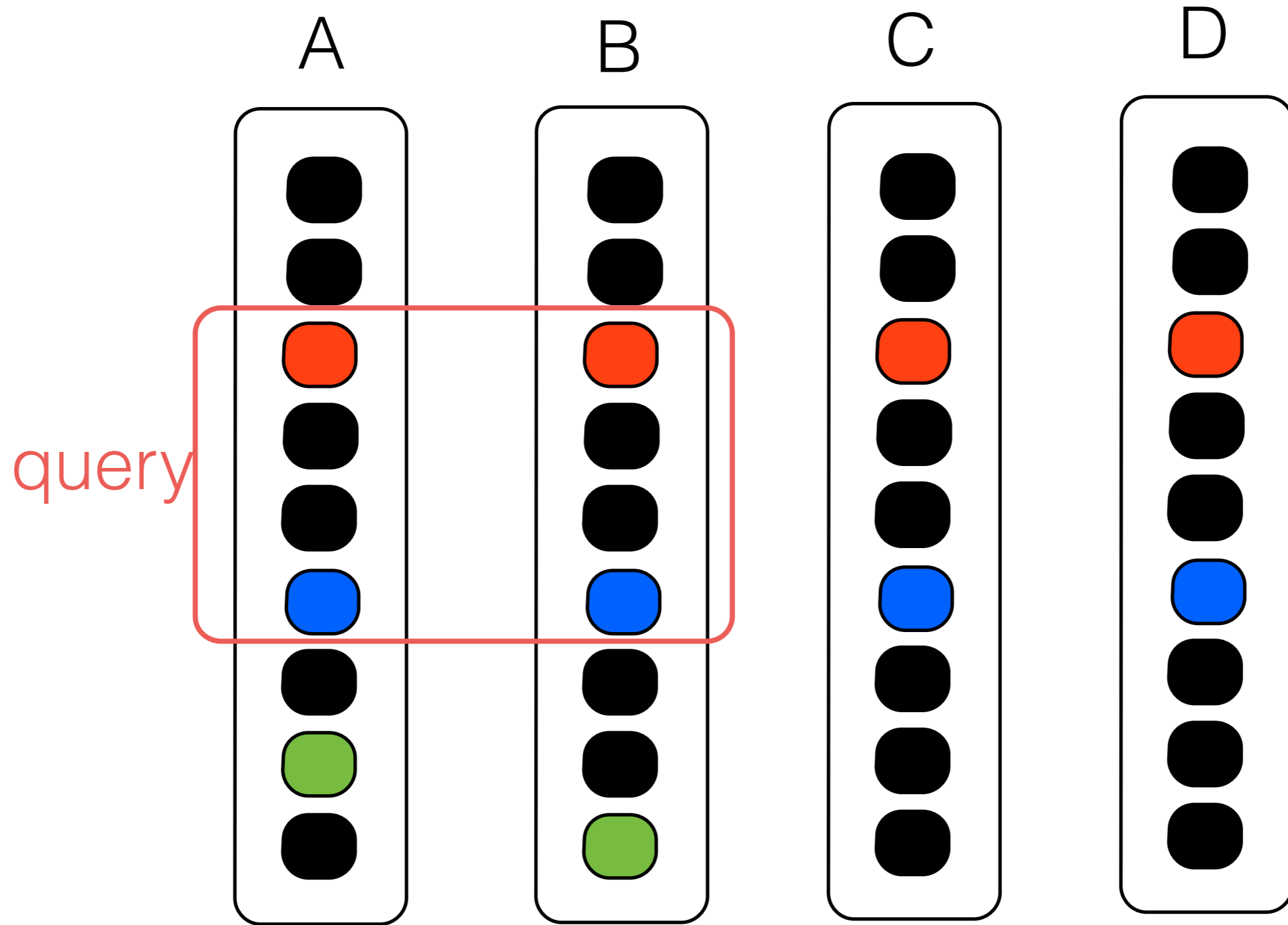
query



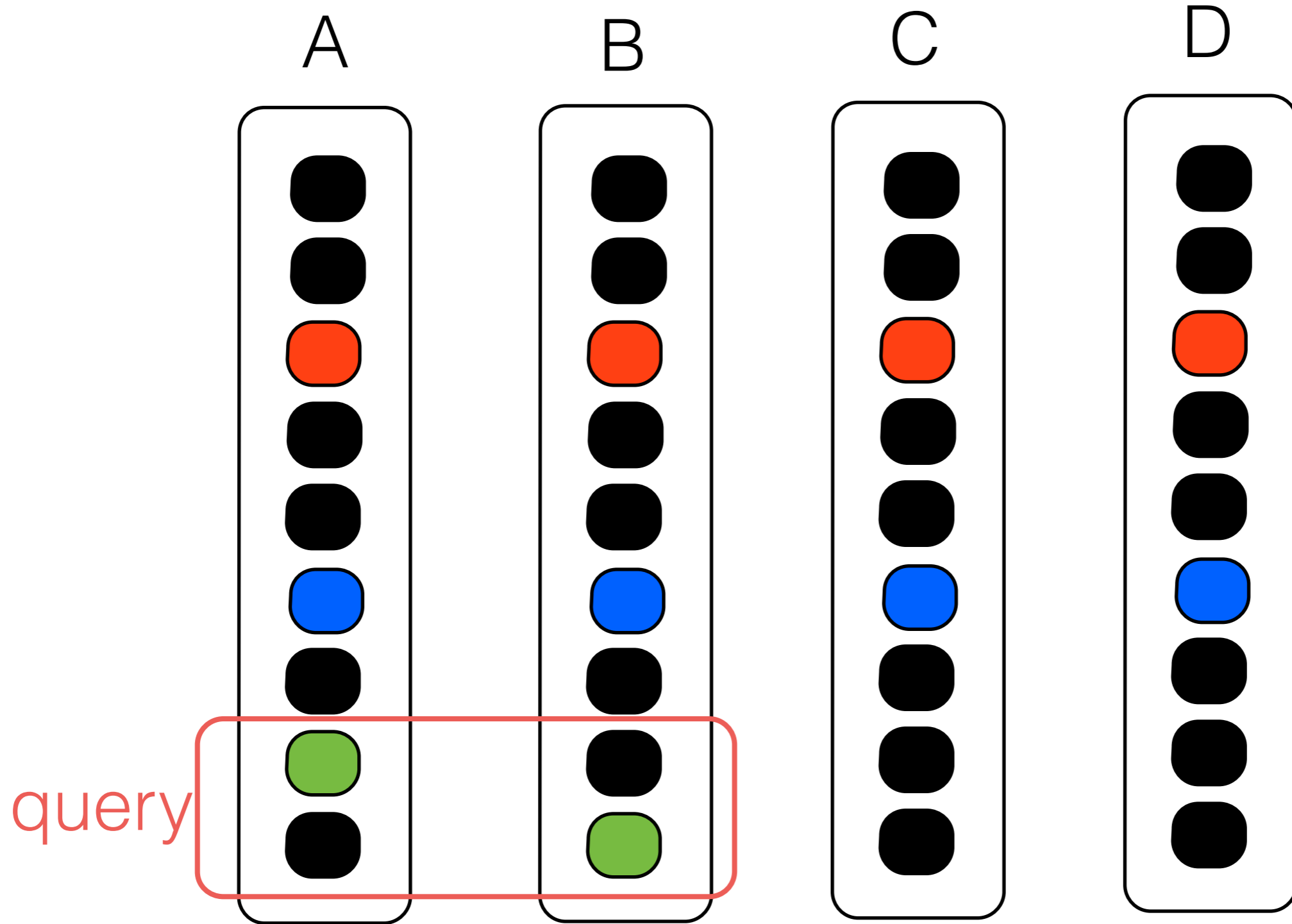
sideways cracking



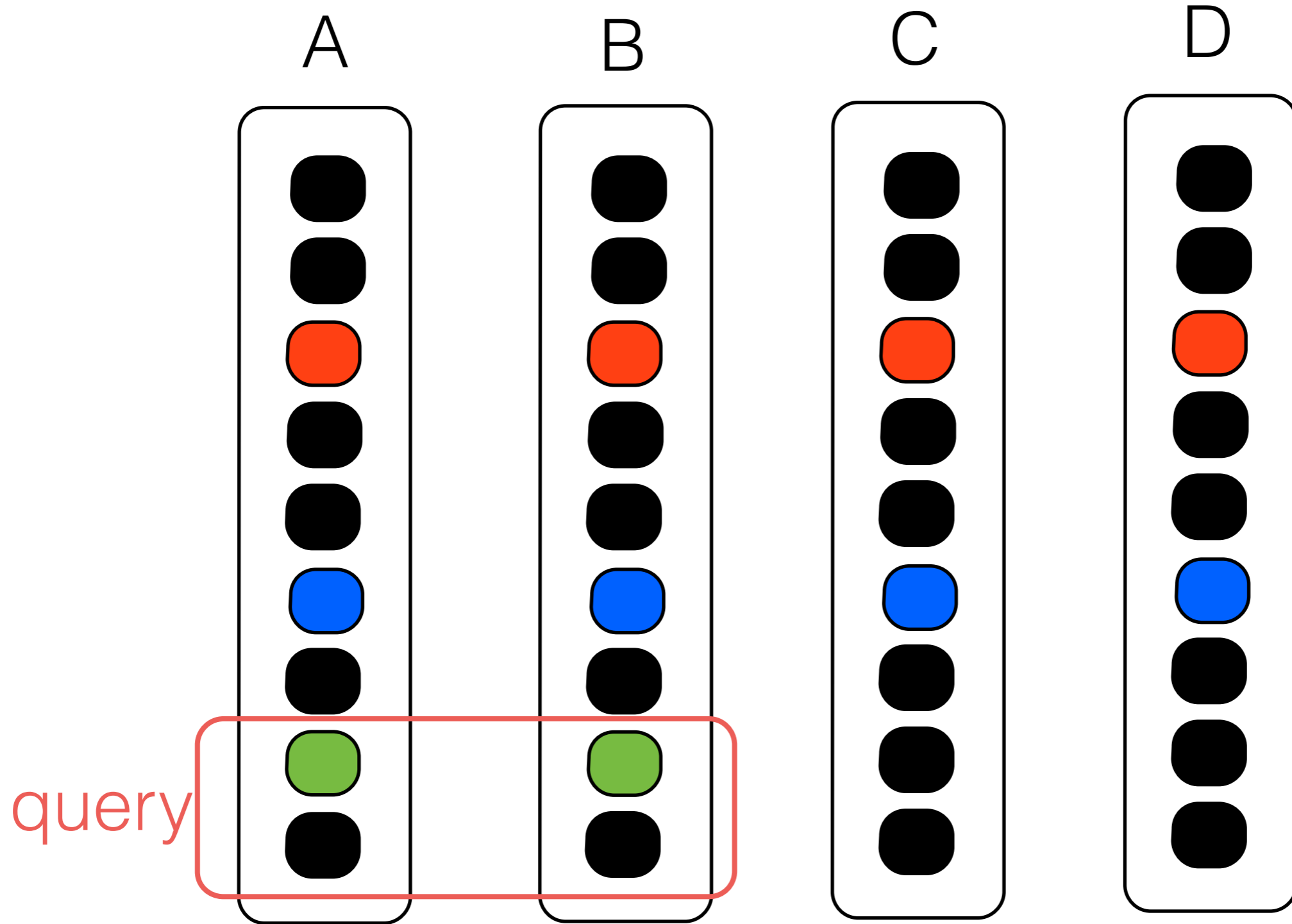
sideways cracking



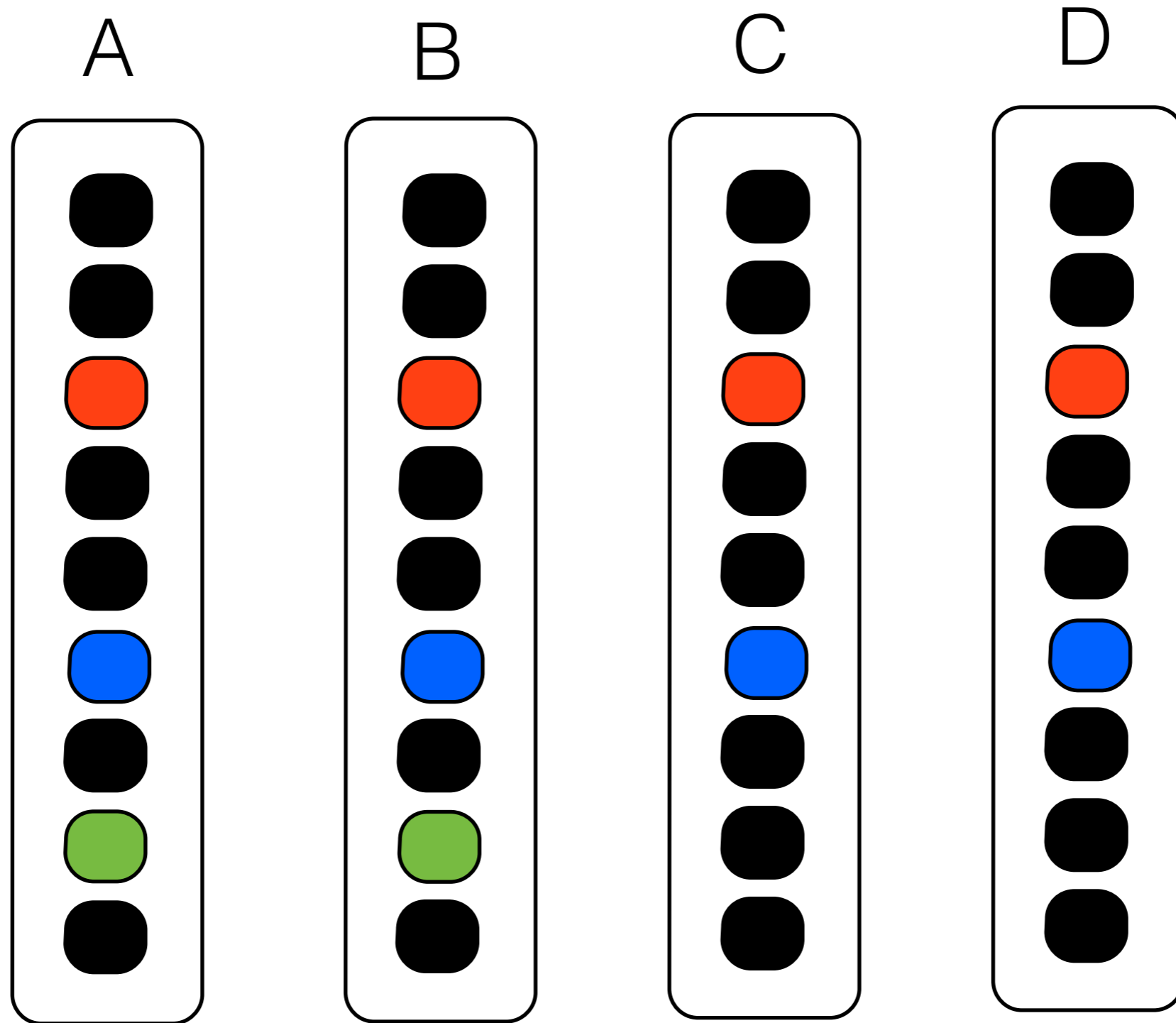
sideways cracking



sideways cracking



sideways cracking



log crack actions and replay to align columns dynamically
replace tuple reconstruction with cracking actions

sideways cracking

Initial state

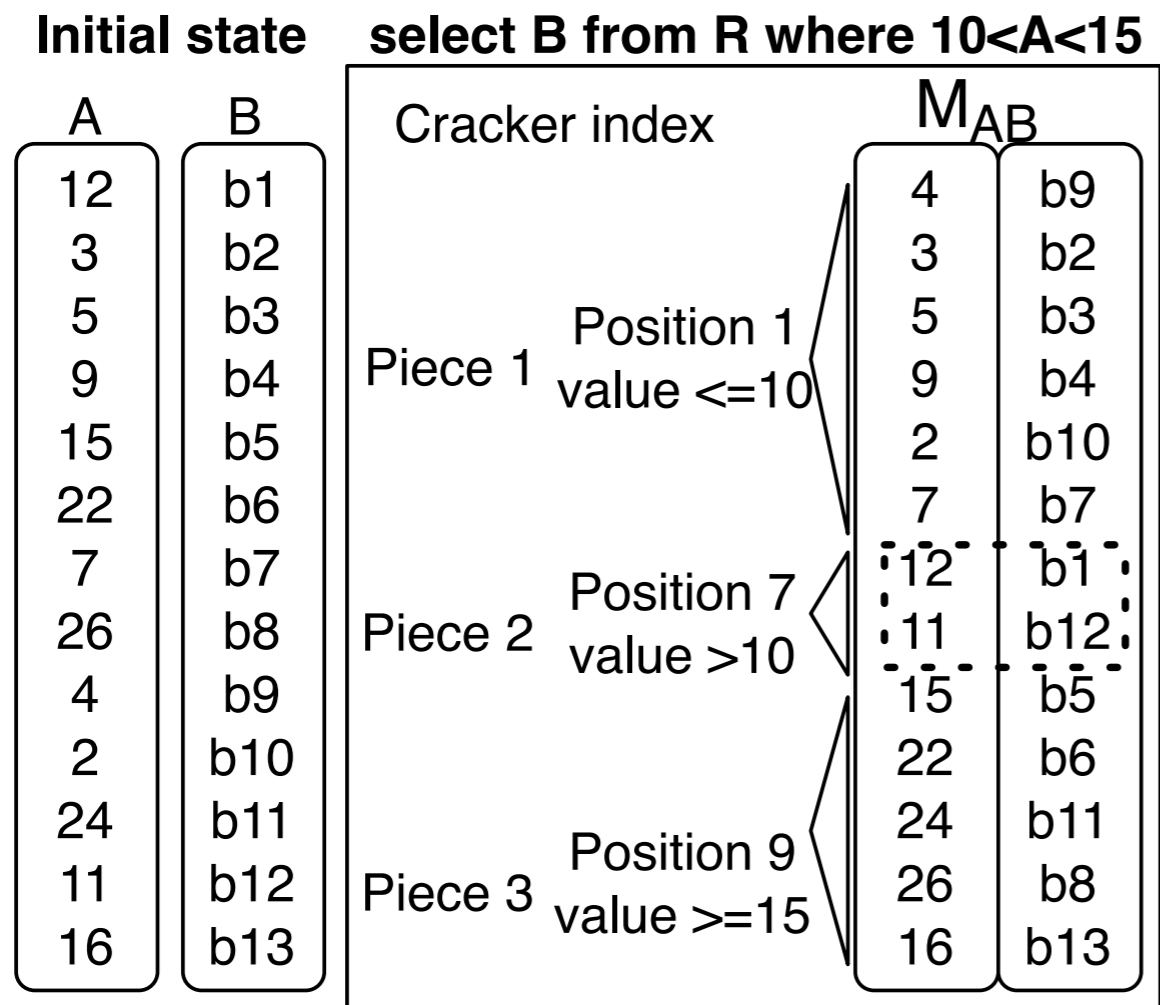
A	B
12	b1
3	b2
5	b3
9	b4
15	b5
22	b6
7	b7
26	b8
4	b9
2	b10
24	b11
11	b12
16	b13

sideways cracking

Initial state **select B from R where $10 < A < 15$**

A	B
12	b1
3	b2
5	b3
9	b4
15	b5
22	b6
7	b7
26	b8
4	b9
2	b10
24	b11
11	b12
16	b13

sideways cracking



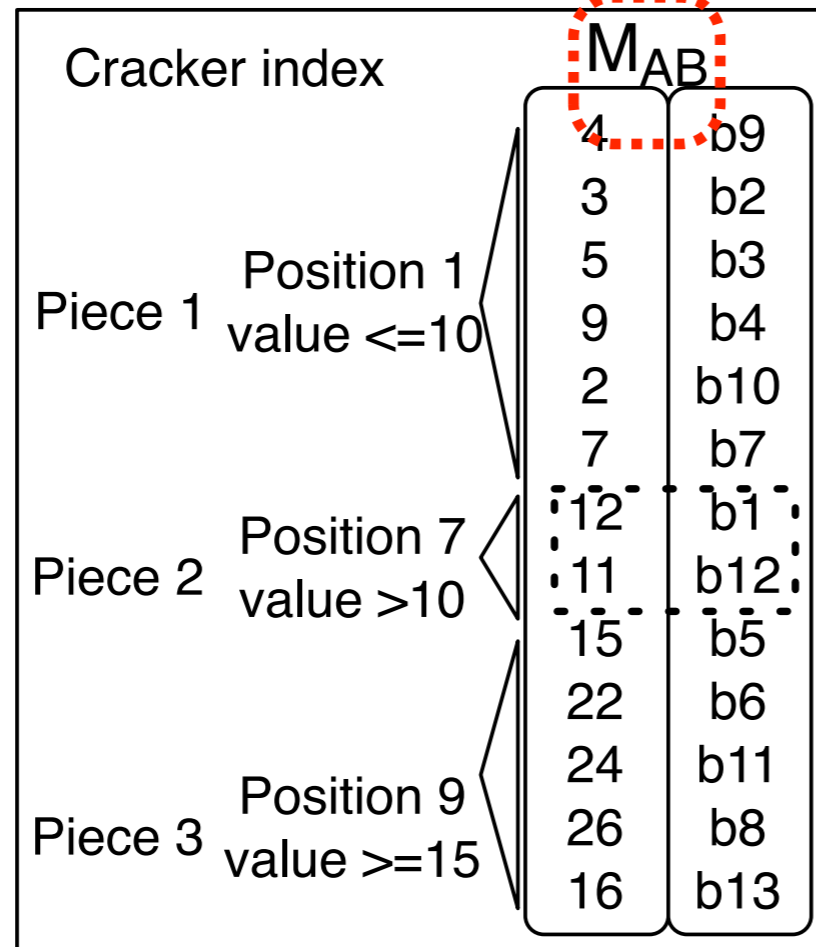
sideways cracking

Cracker Map Head Tail

Initial state

A	B
12	b1
3	b2
5	b3
9	b4
15	b5
22	b6
7	b7
26	b8
4	b9
2	b10
24	b11
11	b12
16	b13

select B from R where $10 < A < 15$



sideways cracking

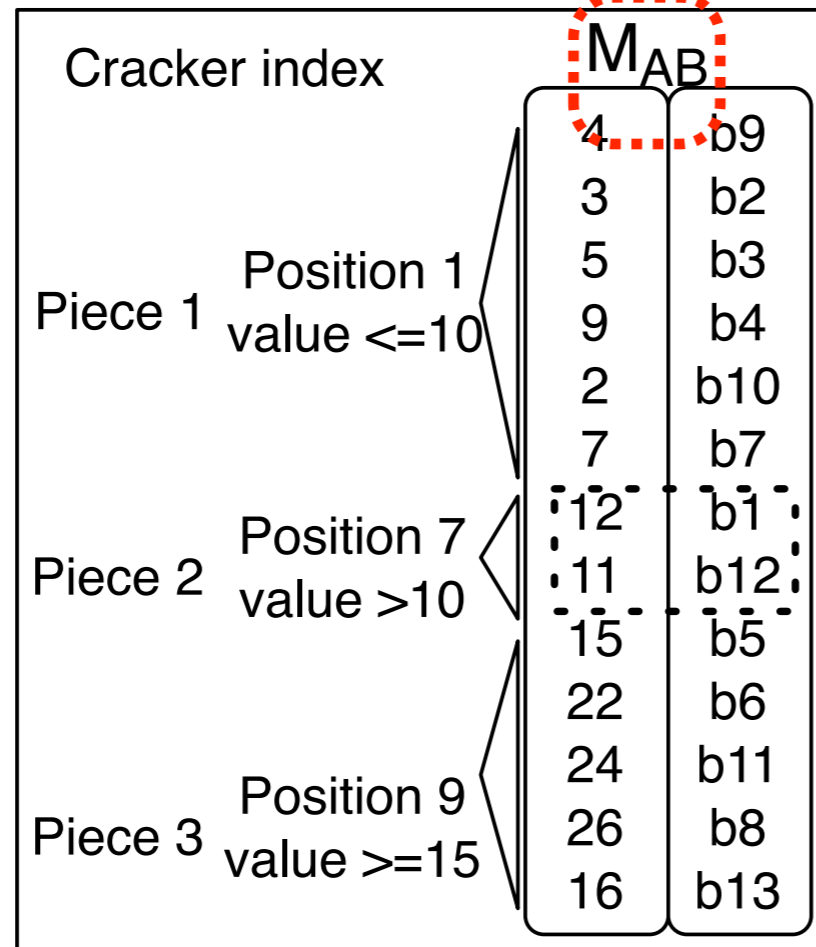
crack based on head, carry tail

Cracker Map
Head Tail

Initial state

A	B
12	b1
3	b2
5	b3
9	b4
15	b5
22	b6
7	b7
26	b8
4	b9
2	b10
24	b11
11	b12
16	b13

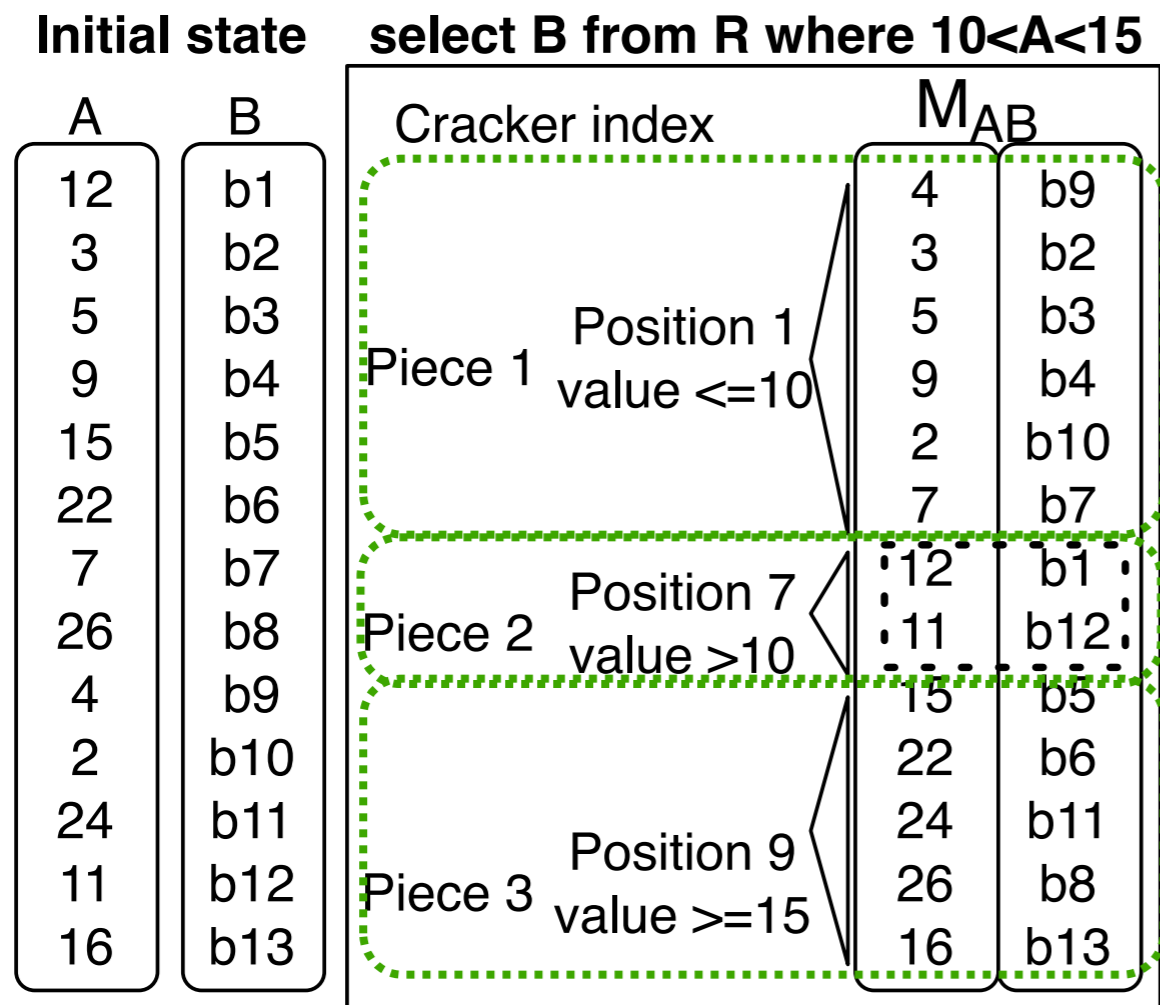
select B from R where $10 < A < 15$



sideways cracking

crack based on head, carry tail

Cracker Map
Head Tail



Cracking knowledge

sideways cracking

crack based on head, carry tail

Cracker Map
Head Tail

Initial state

A	B
12	b1
3	b2
5	b3
9	b4
15	b5
22	b6
7	b7
26	b8
4	b9
2	b10
24	b11
11	b12
16	b13

select B from R where $10 < A < 15$

Cracker index	M_{AB}	
Piece 1 Position 1 value ≤ 10	4	b9
	3	b2
	5	b3
	9	b4
	2	b10
Piece 2 Position 7 value > 10	7	b7
	12	b1
	11	b12
Piece 3 Position 9 value ≥ 15	15	b5
	22	b6
	24	b11
	26	b8
	16	b13

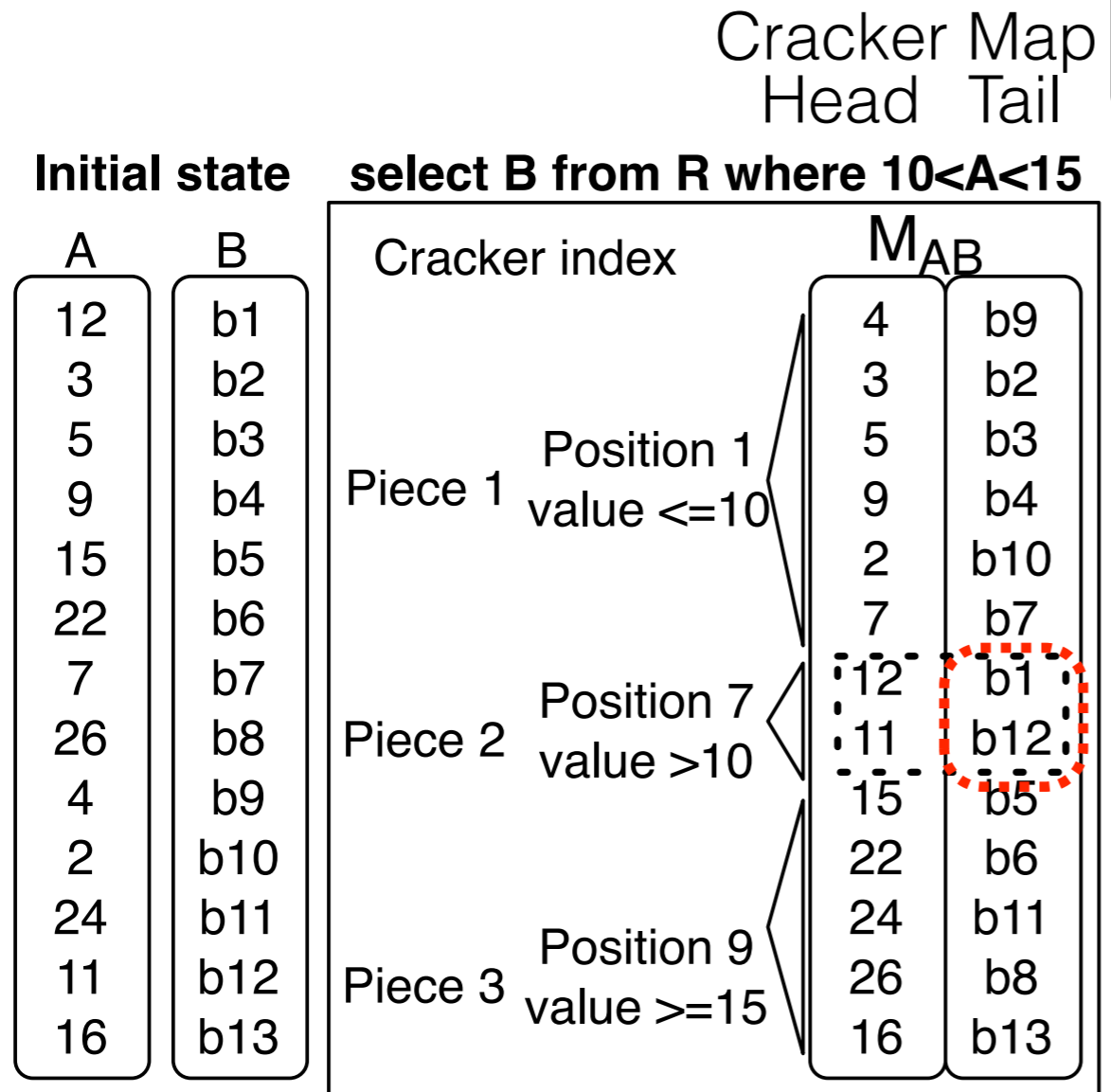


Cracking knowledge

No tuple reconstruction

sideways cracking

crack based on head, carry tail



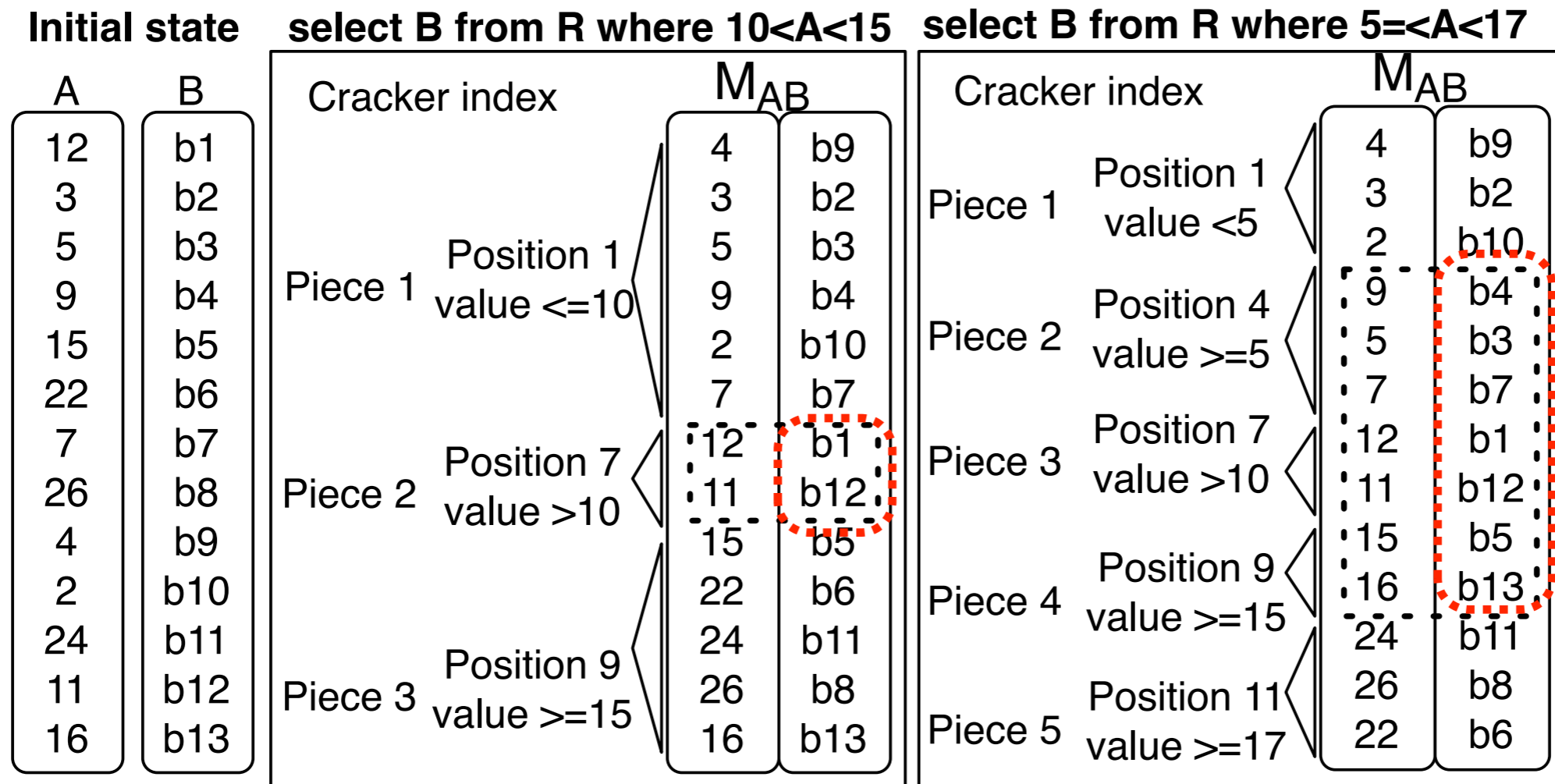
Cracking knowledge No tuple reconstruction

Dynamically/on-the-fly within the select-operator

sideways cracking

Cracker Map
Head Tail

crack based on head, carry tail



Cracking knowledge No tuple reconstruction

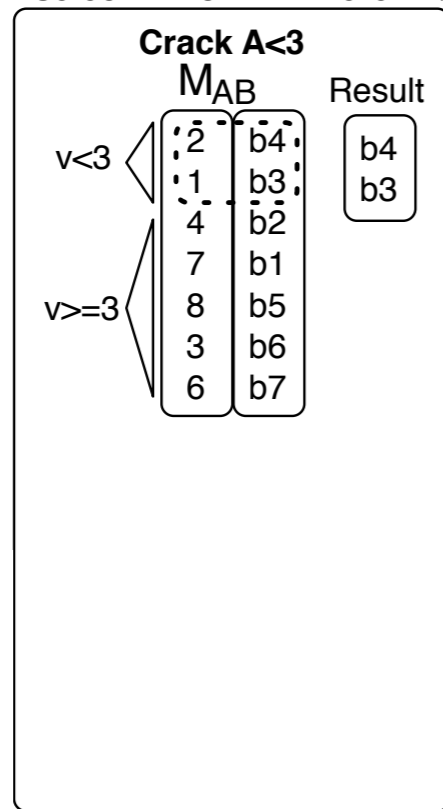
Dynamically/on-the-fly within the select-operator

adaptive alignment

Initial state

A	B	C
7	b1	c1
4	b2	c2
1	b3	c3
2	b4	c4
8	b5	c5
3	b6	c6
6	b7	c7

select B from R where A<3

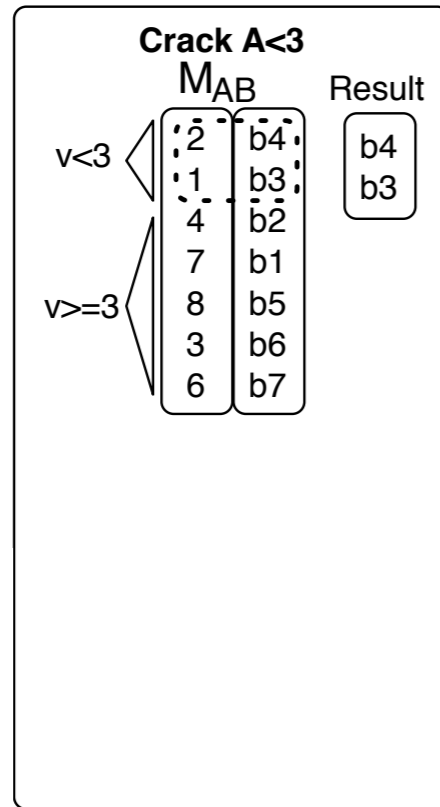


adaptive alignment

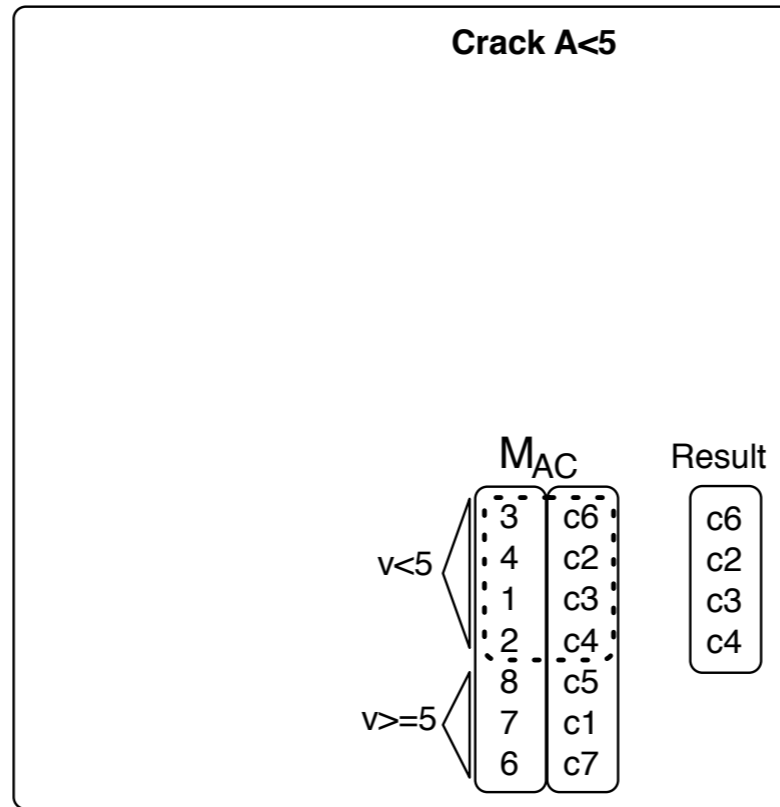
Initial state

A	B	C
7	b1	c1
4	b2	c2
1	b3	c3
2	b4	c4
8	b5	c5
3	b6	c6
6	b7	c7

select B from R where A<3



select C from R where A<5

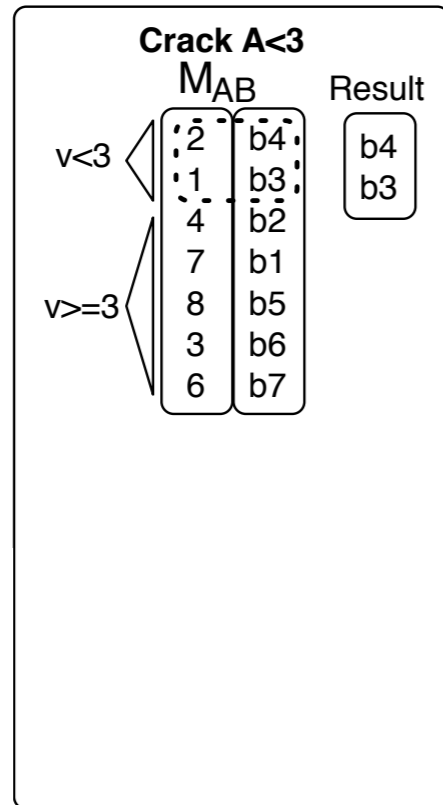


adaptive alignment

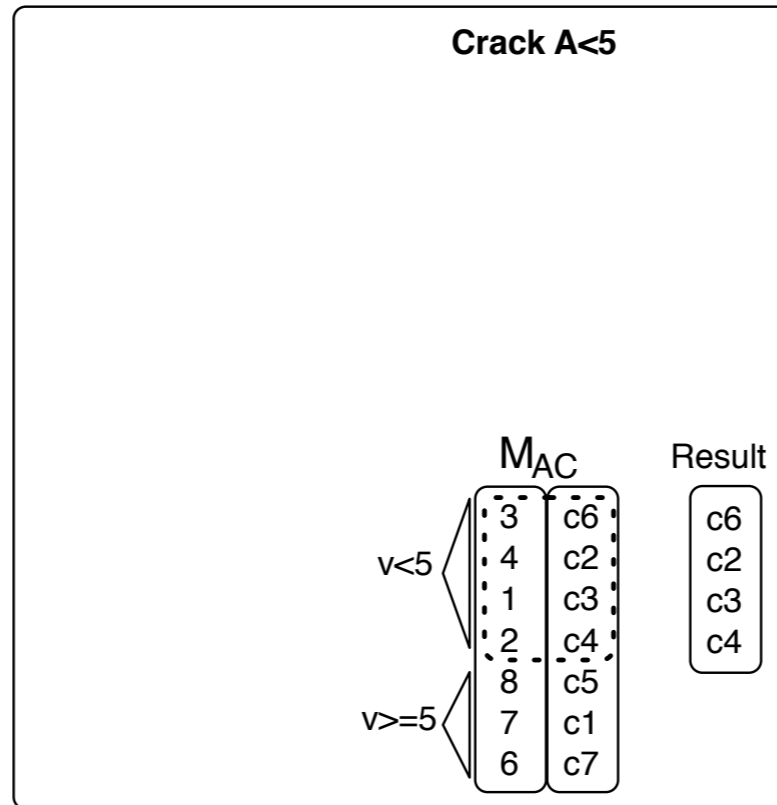
Initial state

A	B	C
7	b1	c1
4	b2	c2
1	b3	c3
2	b4	c4
8	b5	c5
3	b6	c6
6	b7	c7

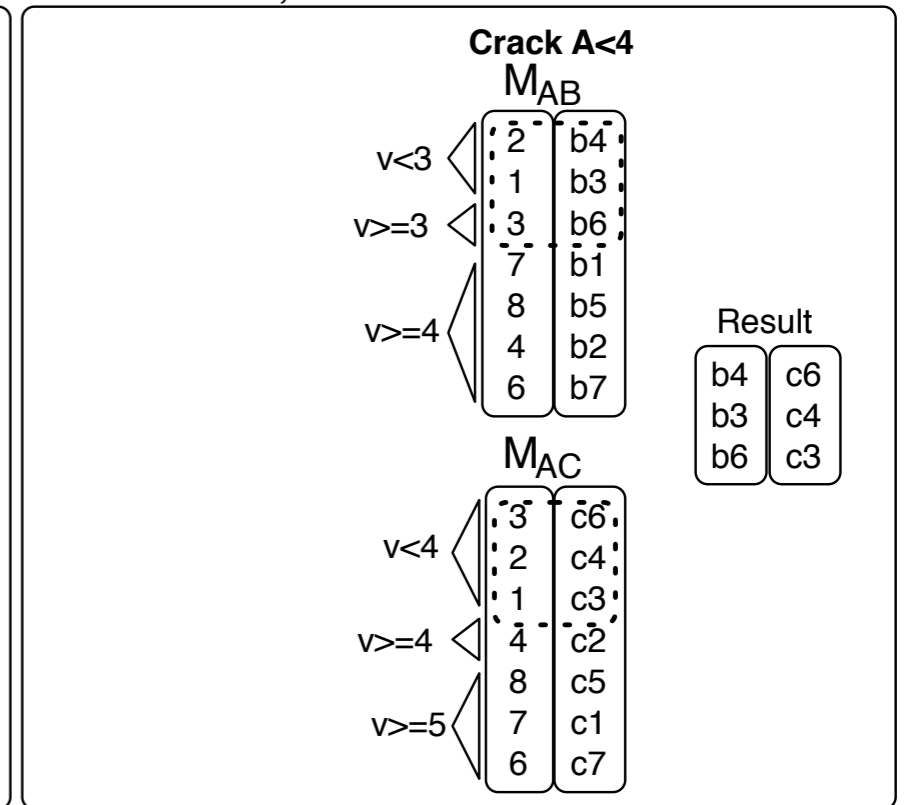
select B from R where A<3



select C from R where A<5



select B,C from R where A<4

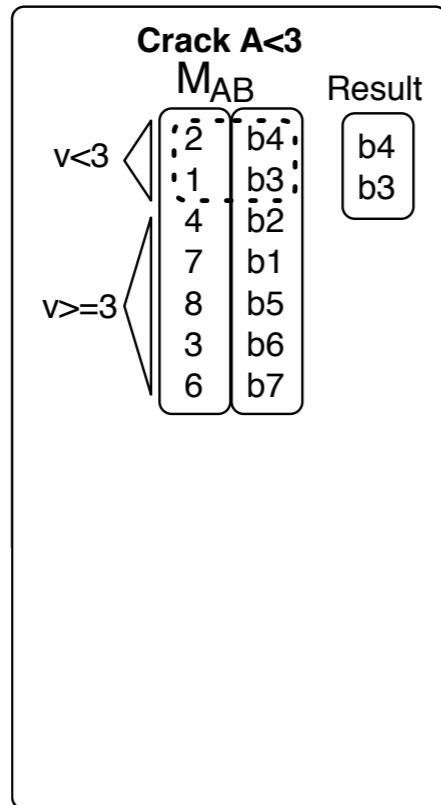


adaptive alignment

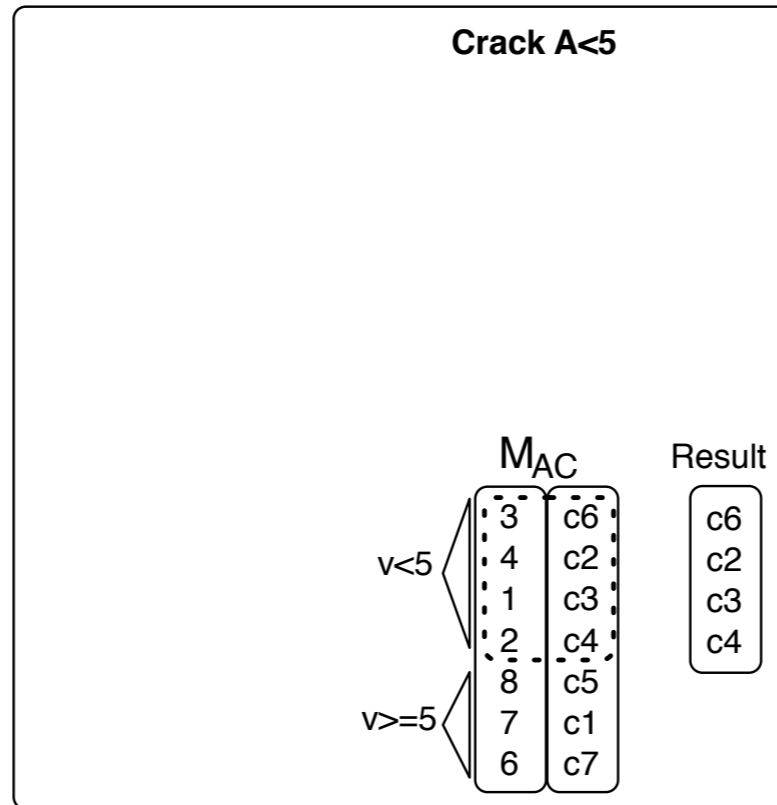
Initial state

A	B	C
7	b1	c1
4	b2	c2
1	b3	c3
2	b4	c4
8	b5	c5
3	b6	c6
6	b7	c7

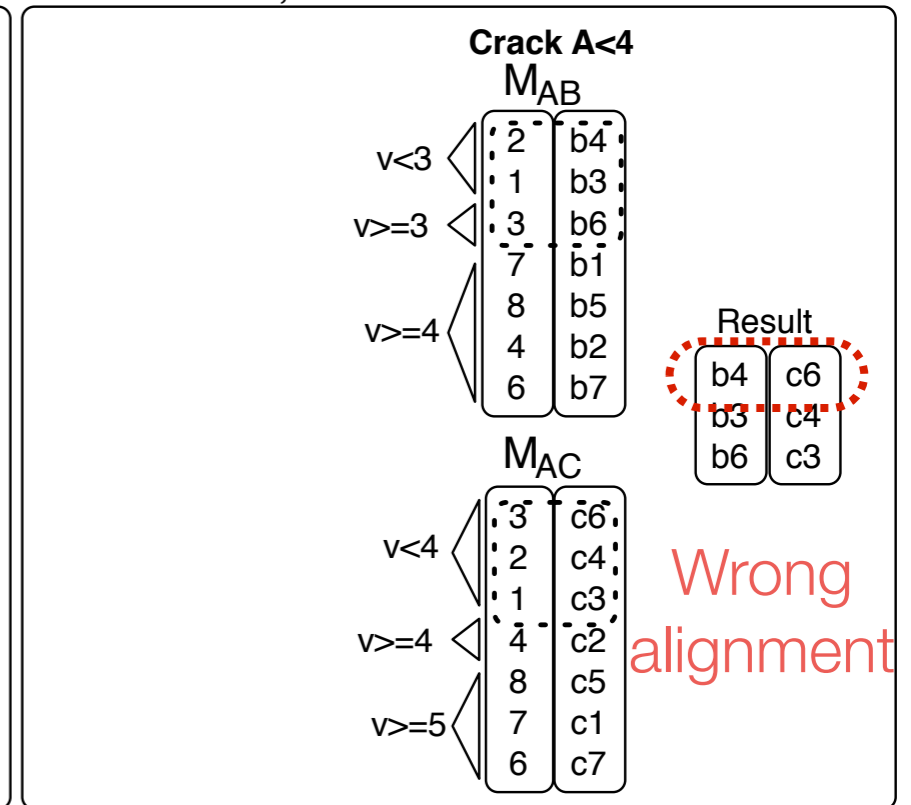
select B from R where A<3



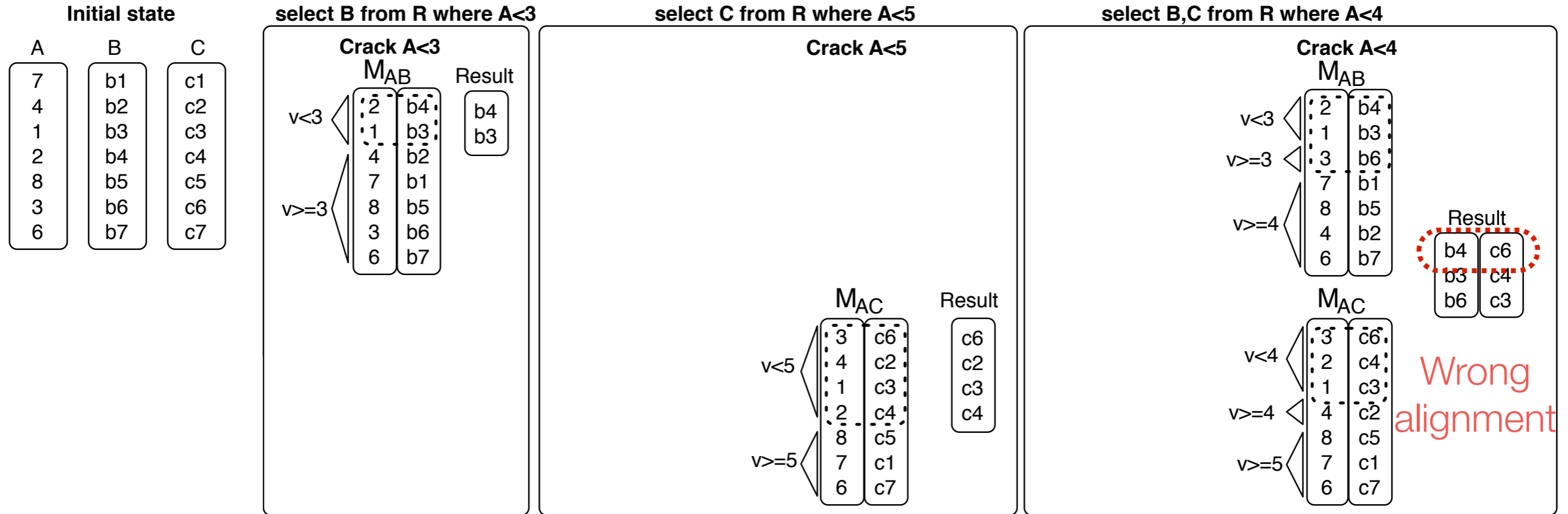
select C from R where A<5



select B,C from R where A<4

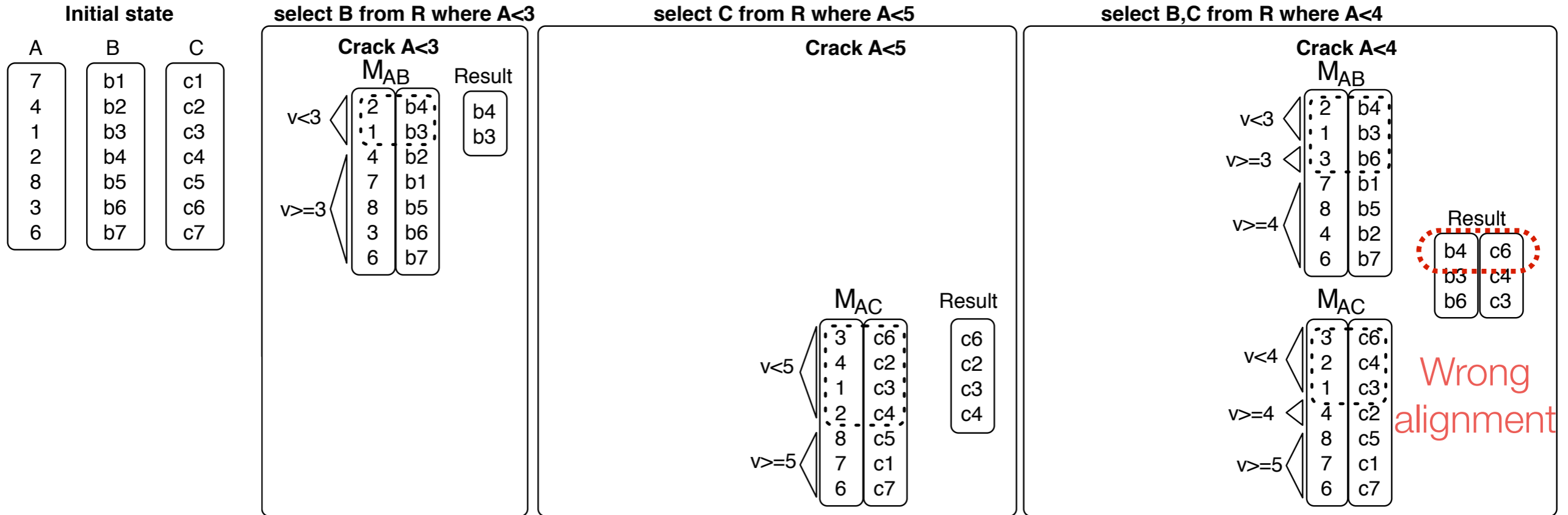


adaptive alignment



perform the same cracks and in the same order on all maps with the same head

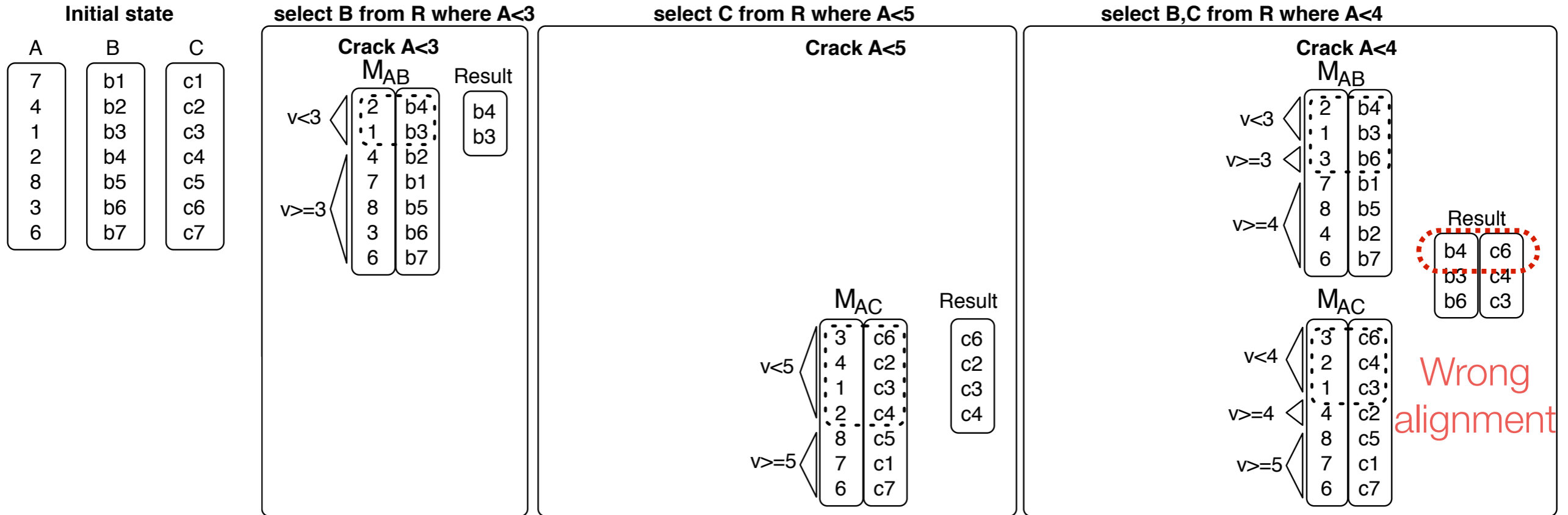
adaptive alignment



perform the same cracks and in the same order on all maps with the same head

on-line alignment
touch/load everything, always

adaptive alignment

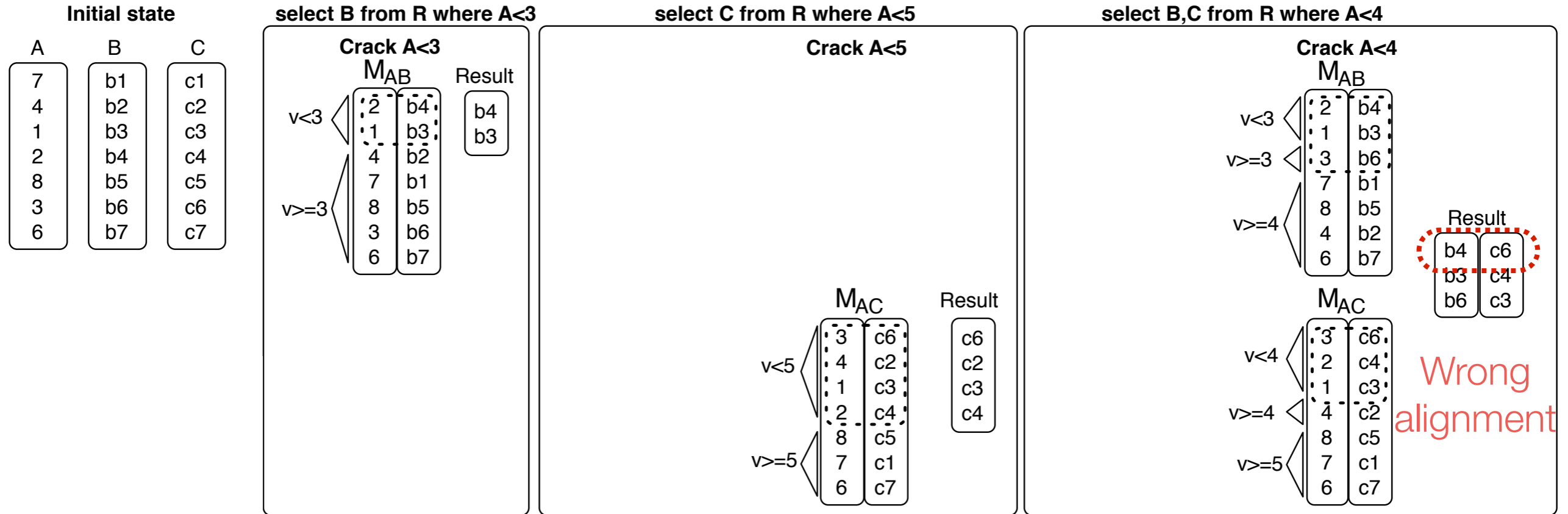


perform the same cracks and in the same order on all maps with the same head

on-line alignment
touch/load everything, always



adaptive alignment



perform the same cracks and in the same order on all maps with the same head

on-line alignment
touch/load everything, always

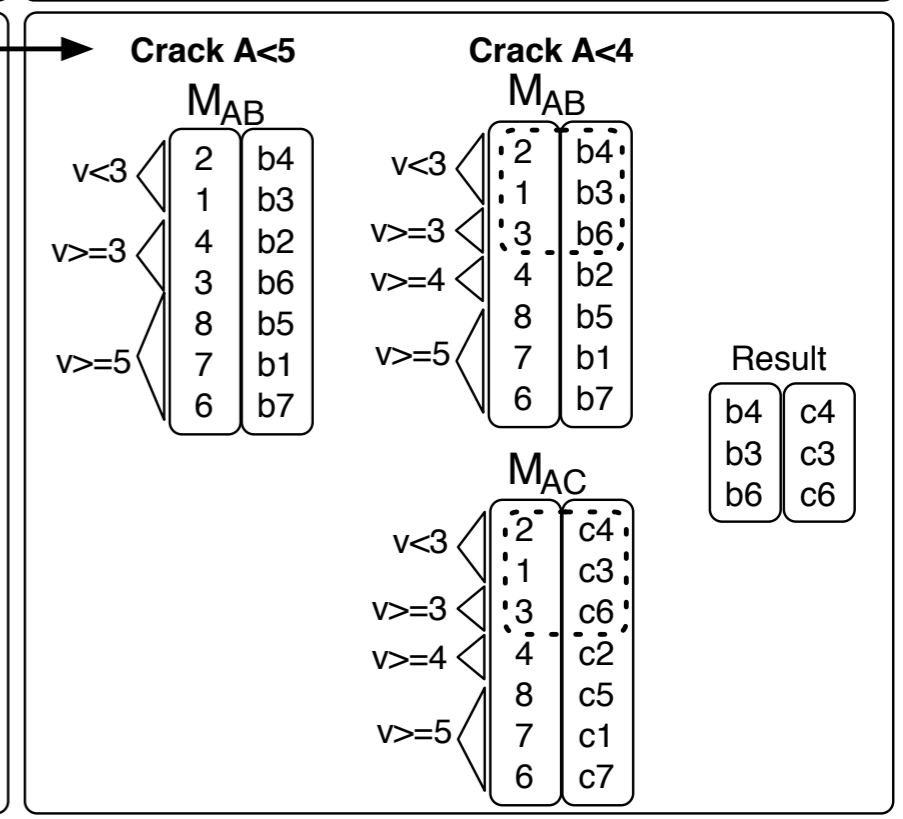
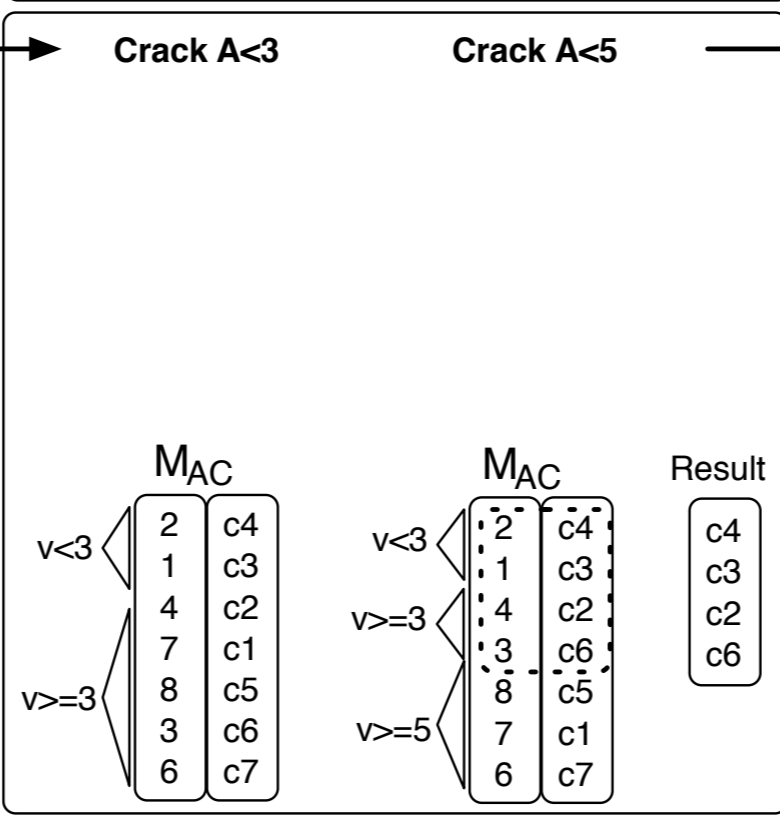
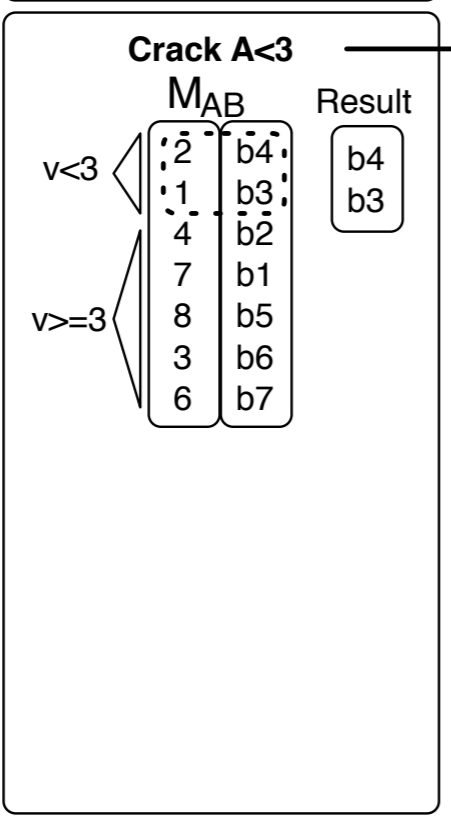
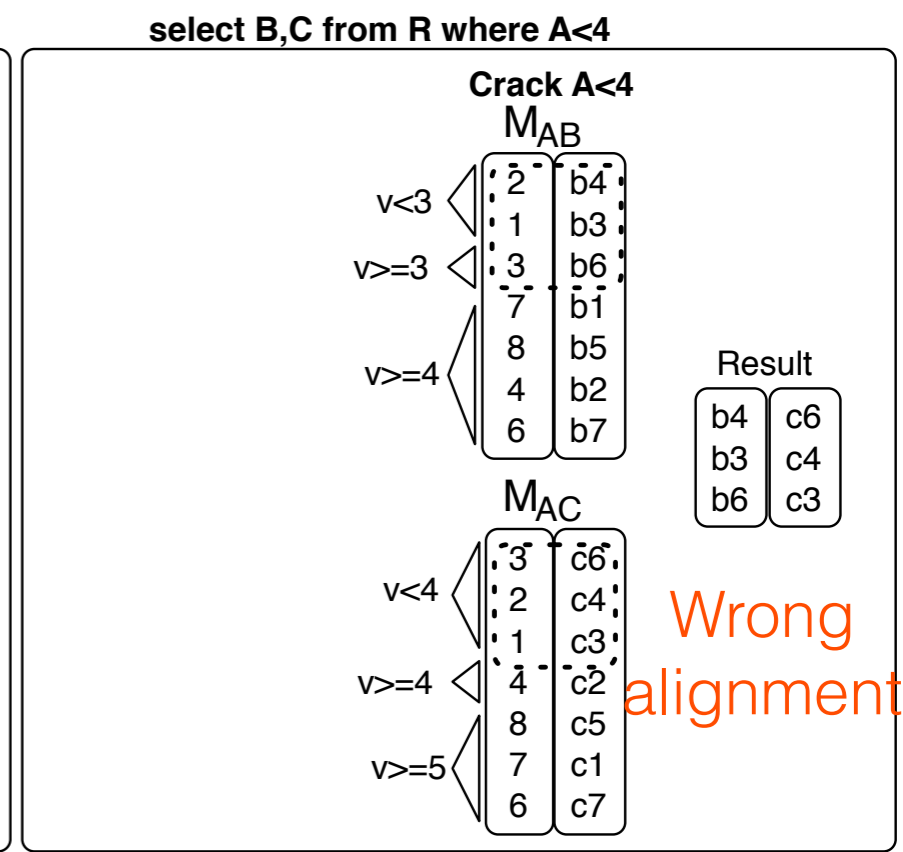
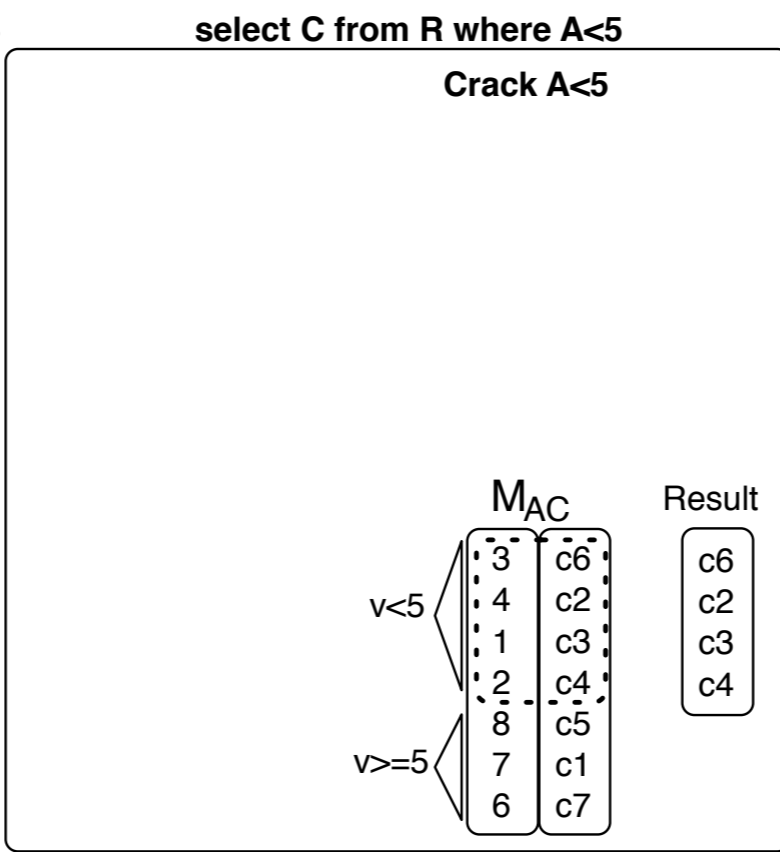
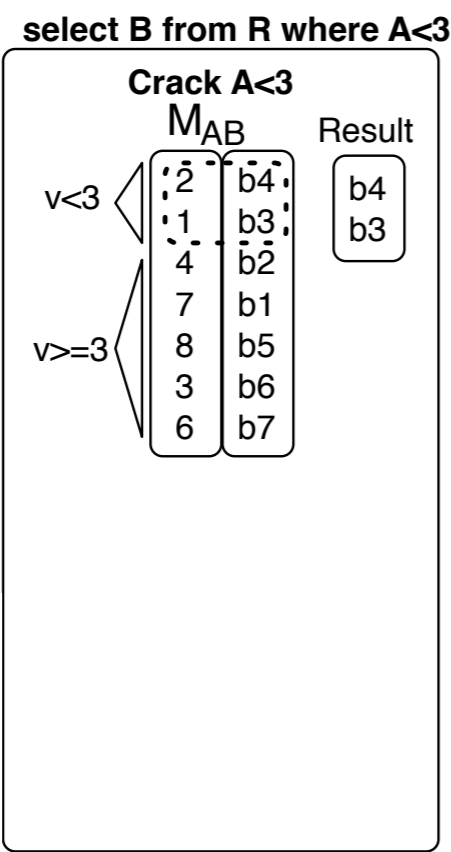


remember and replay cracks across columns

adaptive alignment

Initial state

A	B	C
7	b1	c1
4	b2	c2
1	b3	c3
2	b4	c4
8	b5	c5
3	b6	c6
6	b7	c7

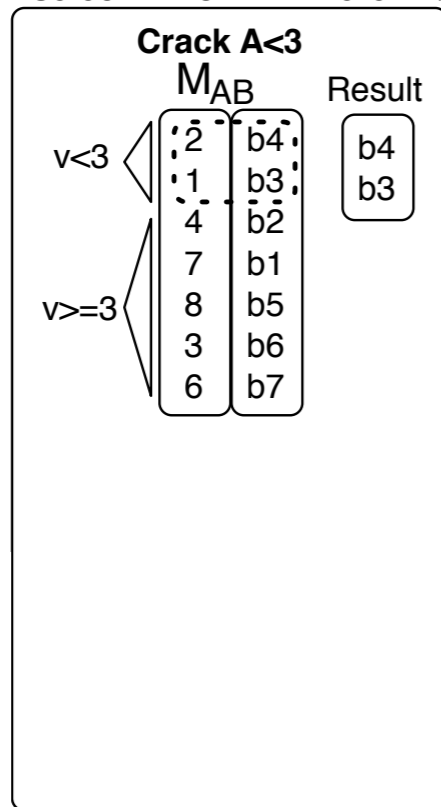


adaptive alignment

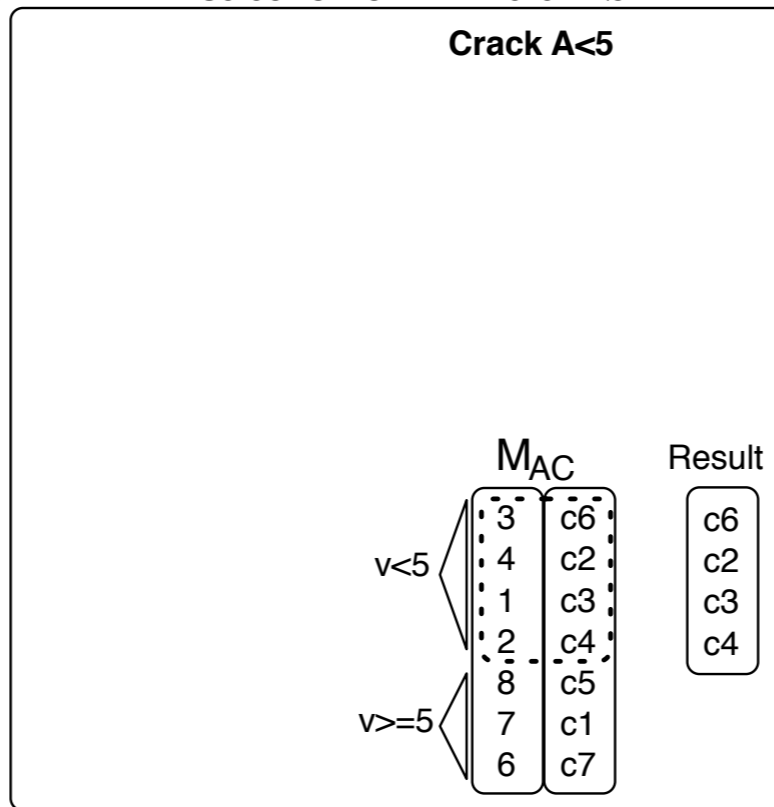
Initial state

A	B	C
7	b1	c1
4	b2	c2
1	b3	c3
2	b4	c4
8	b5	c5
3	b6	c6
6	b7	c7

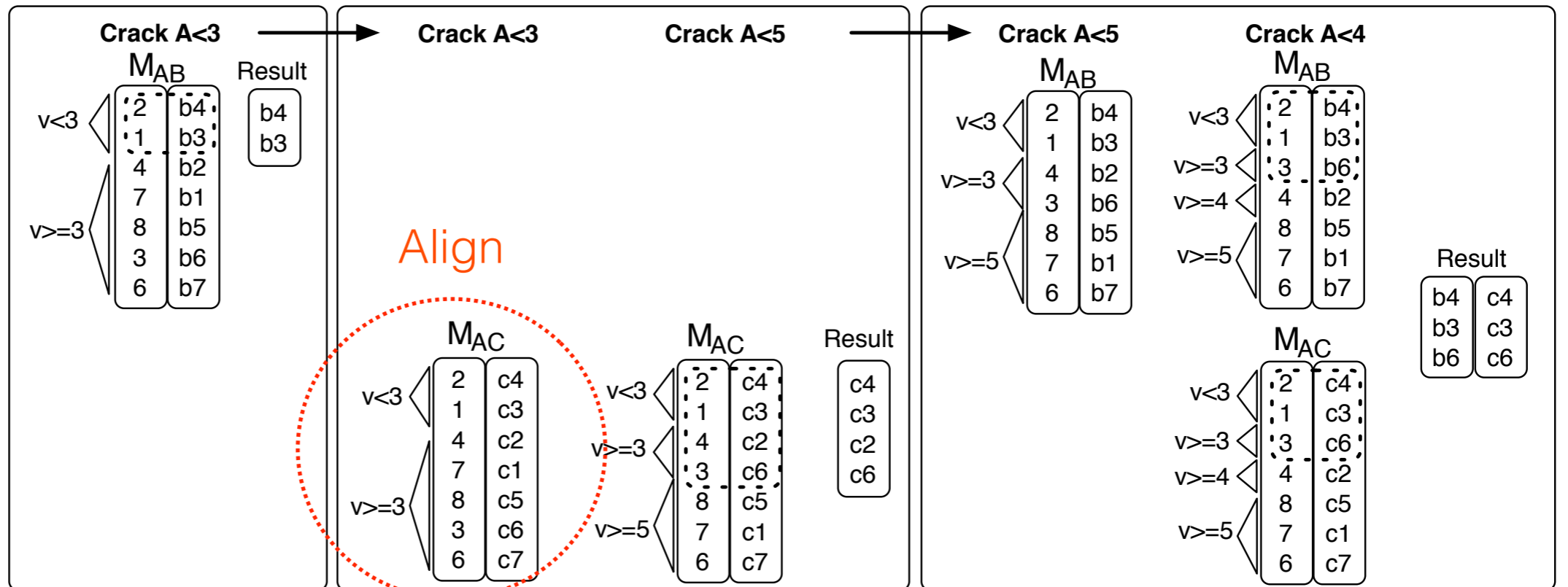
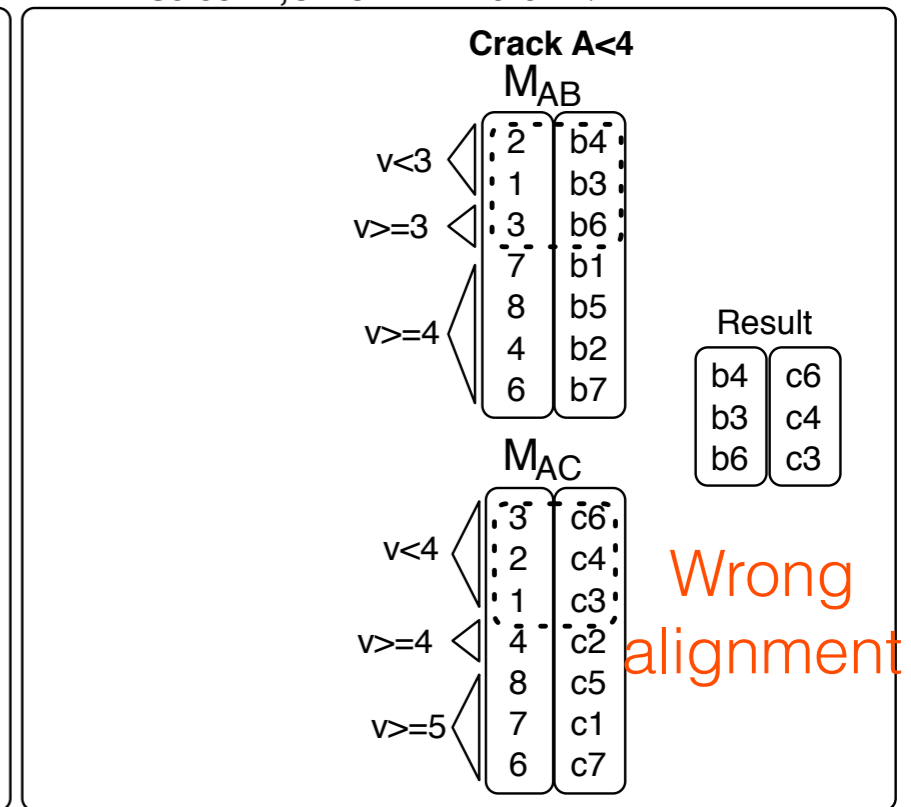
select B from R where A<3



select C from R where A<5



select B,C from R where A<4

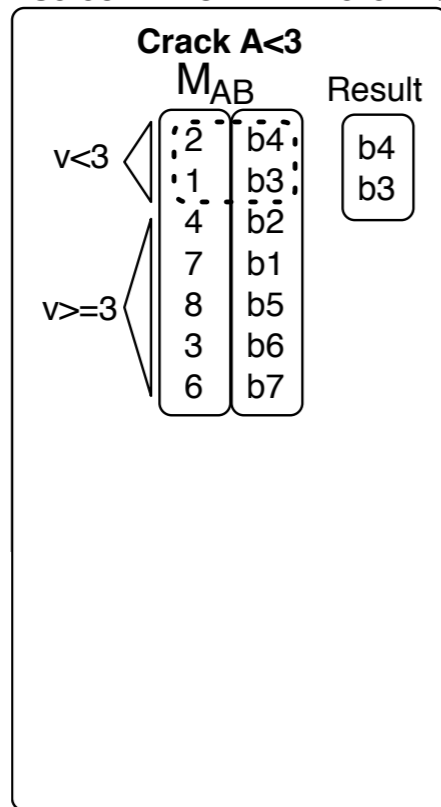


adaptive alignment

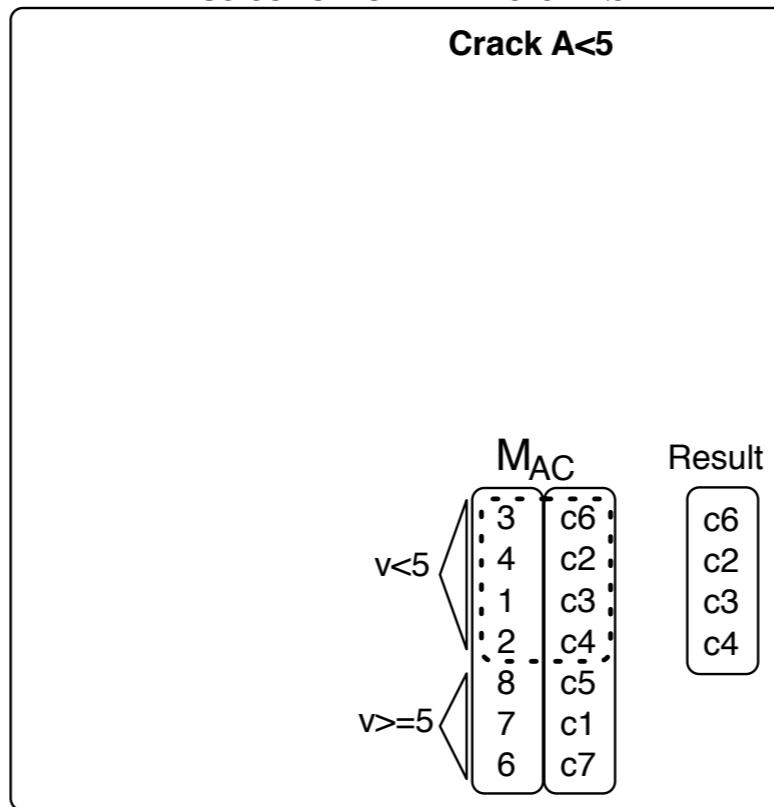
Initial state

A	B	C
7	b1	c1
4	b2	c2
1	b3	c3
2	b4	c4
8	b5	c5
3	b6	c6
6	b7	c7

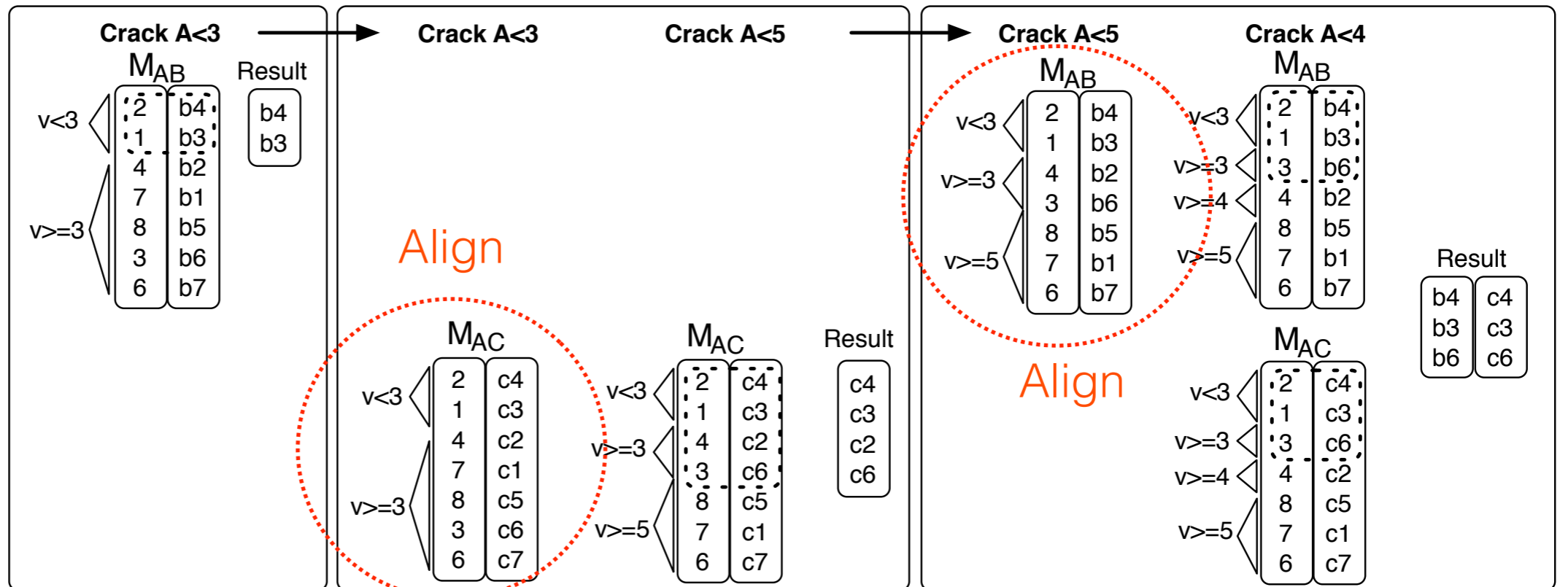
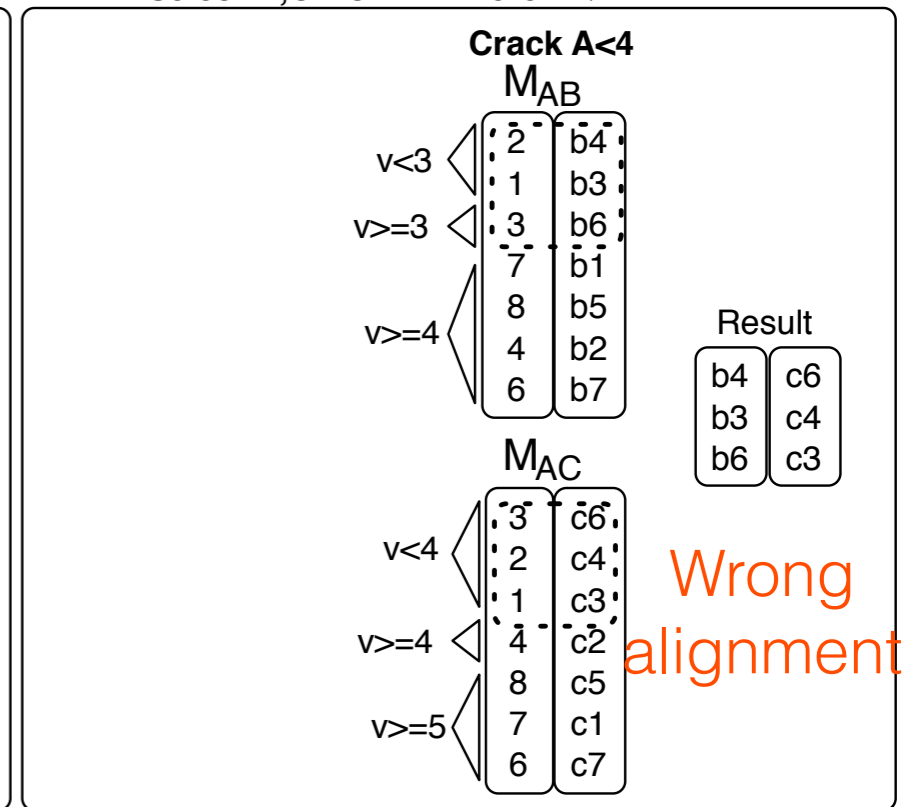
select B from R where A<3



select C from R where A<5



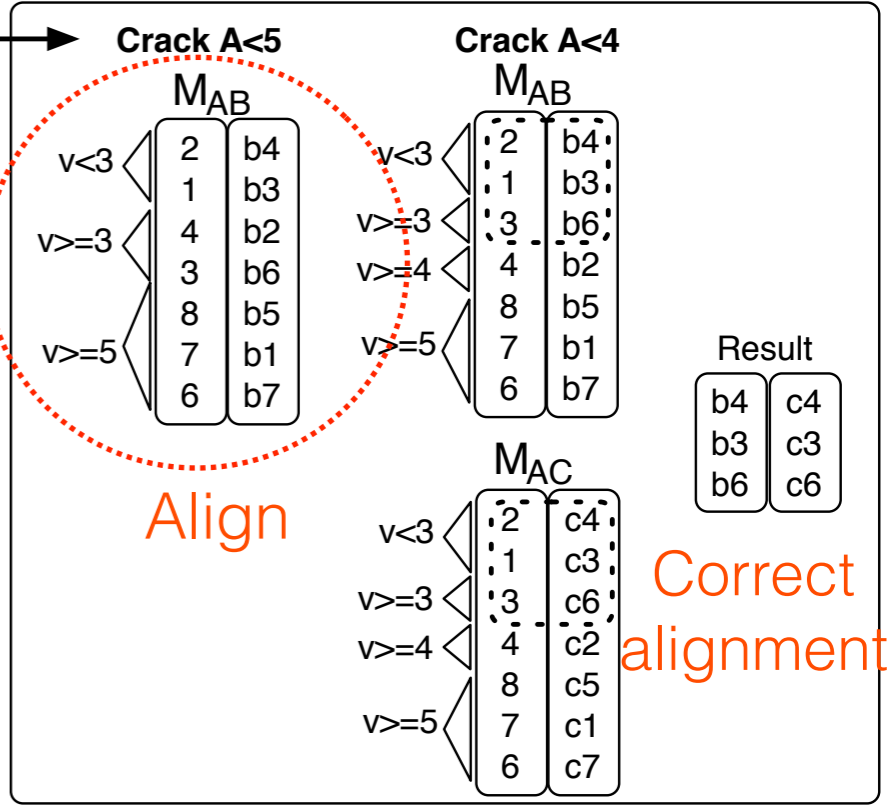
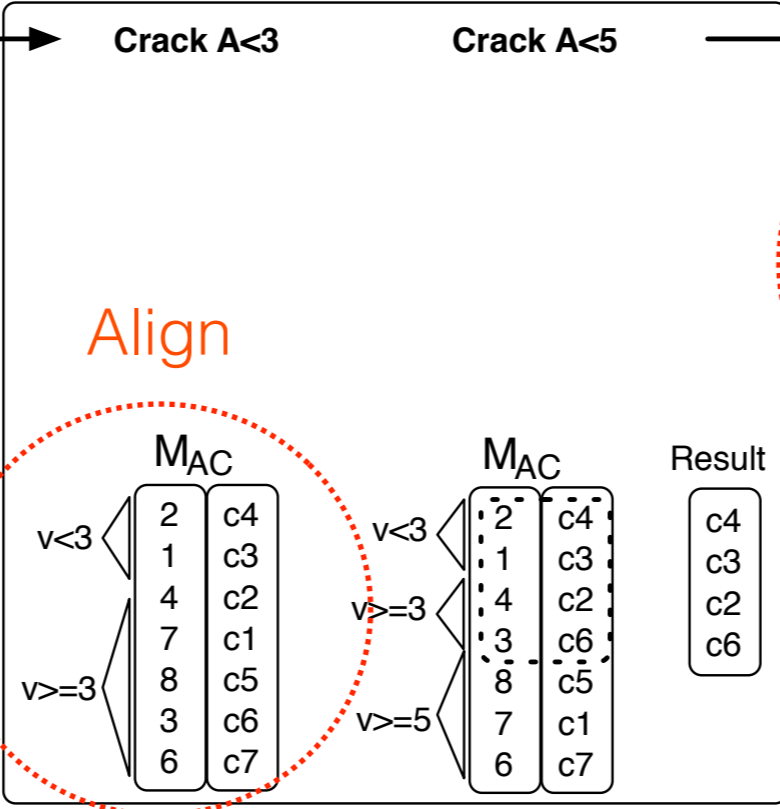
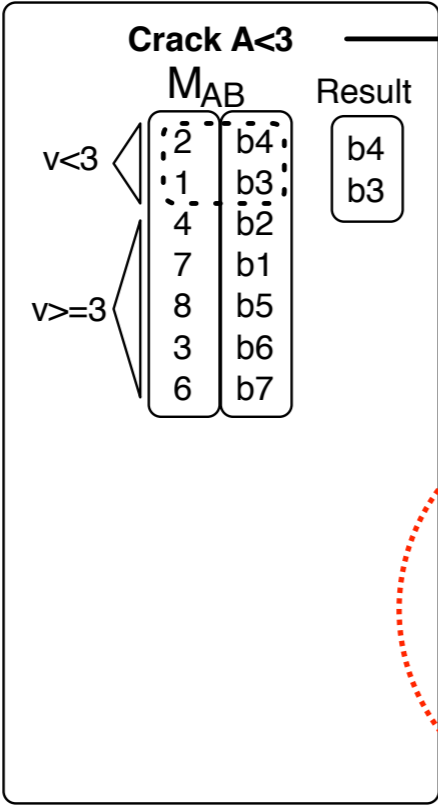
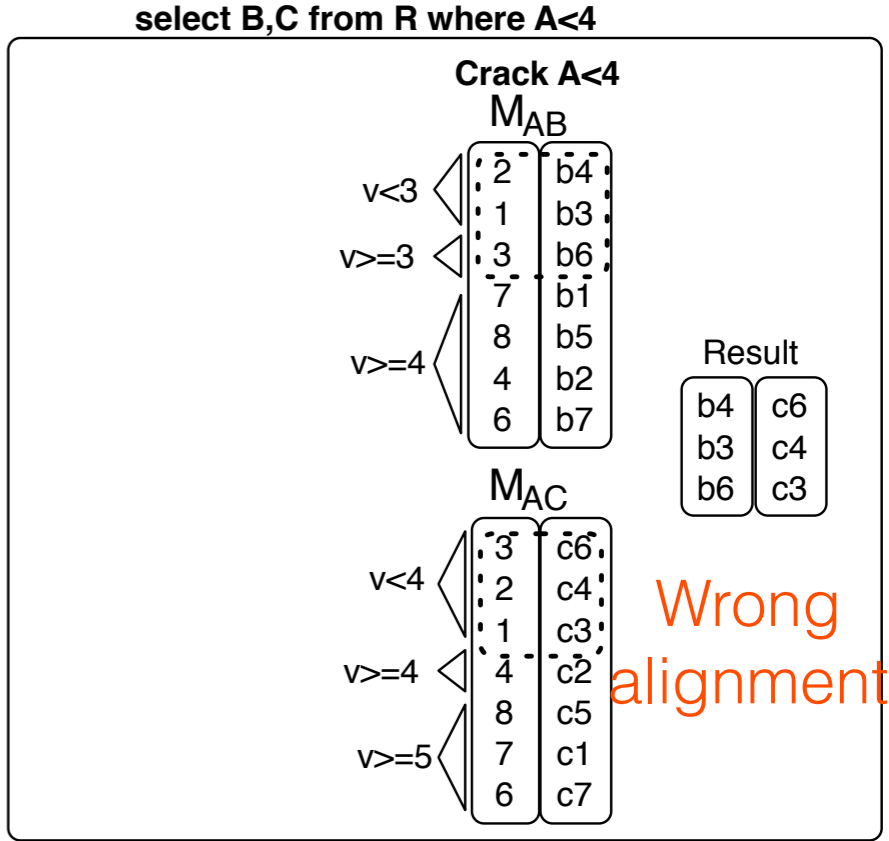
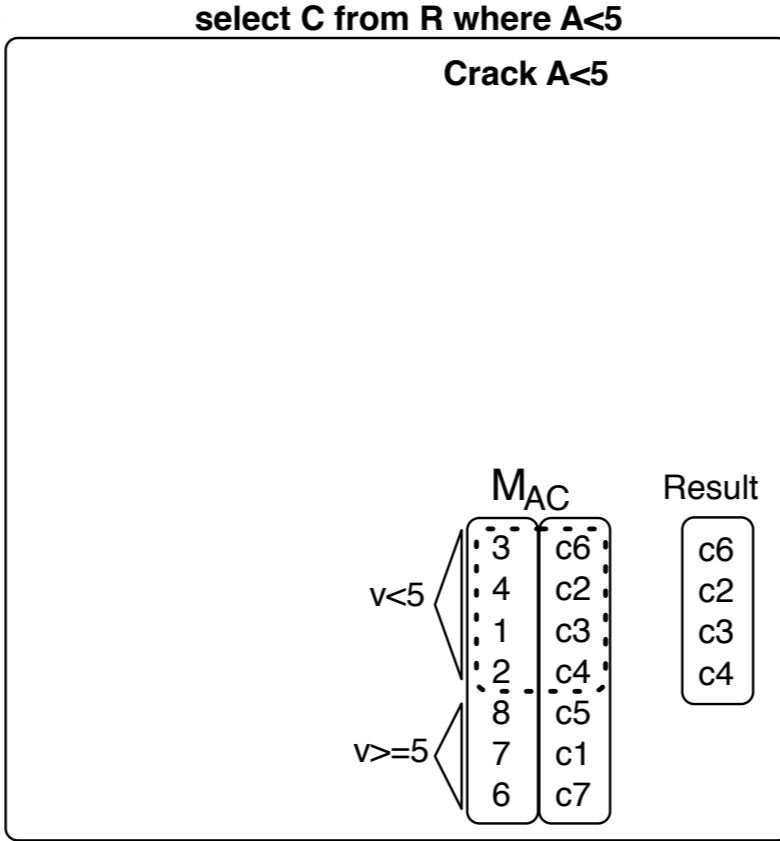
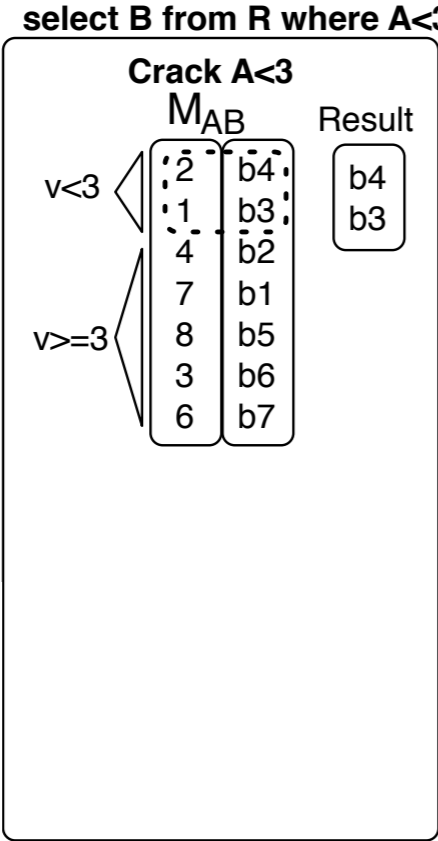
select B,C from R where A<4



adaptive alignment

Initial state

A	B	C
7	b1	c1
4	b2	c2
1	b3	c3
2	b4	c4
8	b5	c5
3	b6	c6
6	b7	c7

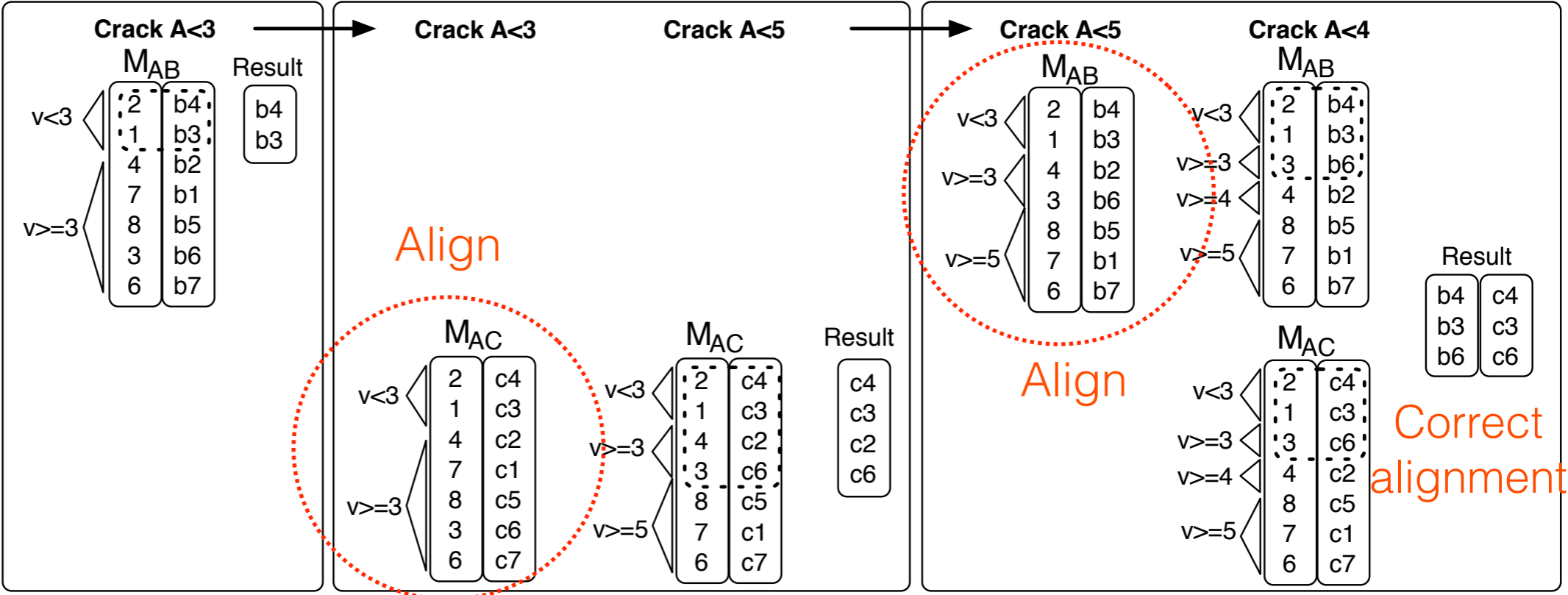
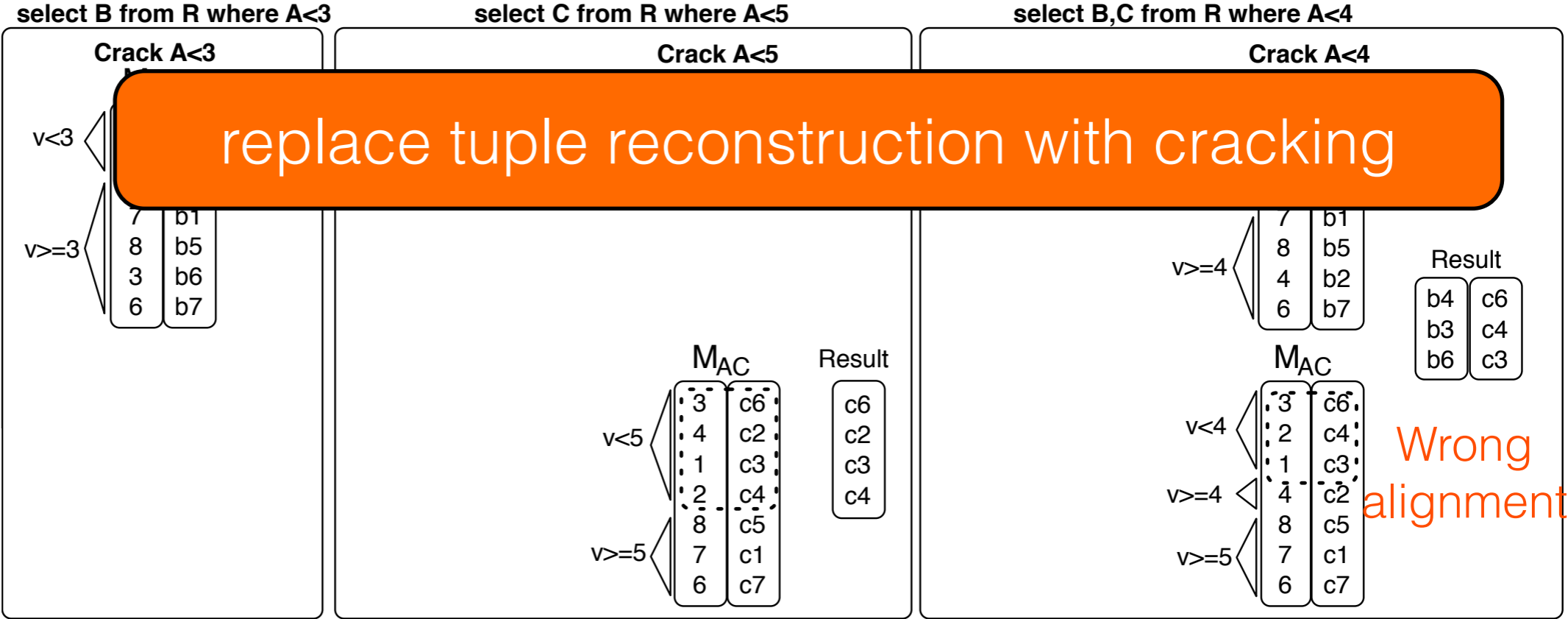


adaptive alignment

Initial state

A	B	C
7	b1	c1
4	b2	c2
1	b3	c3
2	b4	c4
8	b5	c5
3	b6	c6
6	b7	c7

replace tuple reconstruction with cracking



multi-selections

wider maps...but too many combinations

maps of different maps sets lead to alignment problems

select D from R where $3 < A < 10$ and $4 < B < 8$ and $1 < C < 7$

multi-selections

wider maps...but too many combinations

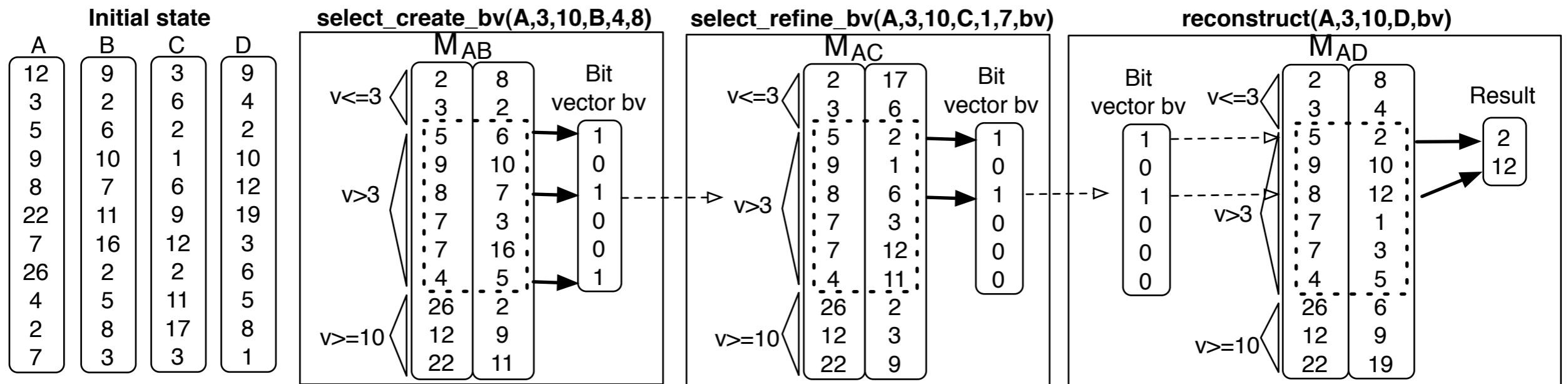
maps of different maps sets lead to alignment problems

use a single map set and exploit bit-vectors

select D from R where $3 < A < 10$ and $4 < B < 8$ and $1 < C < 7$

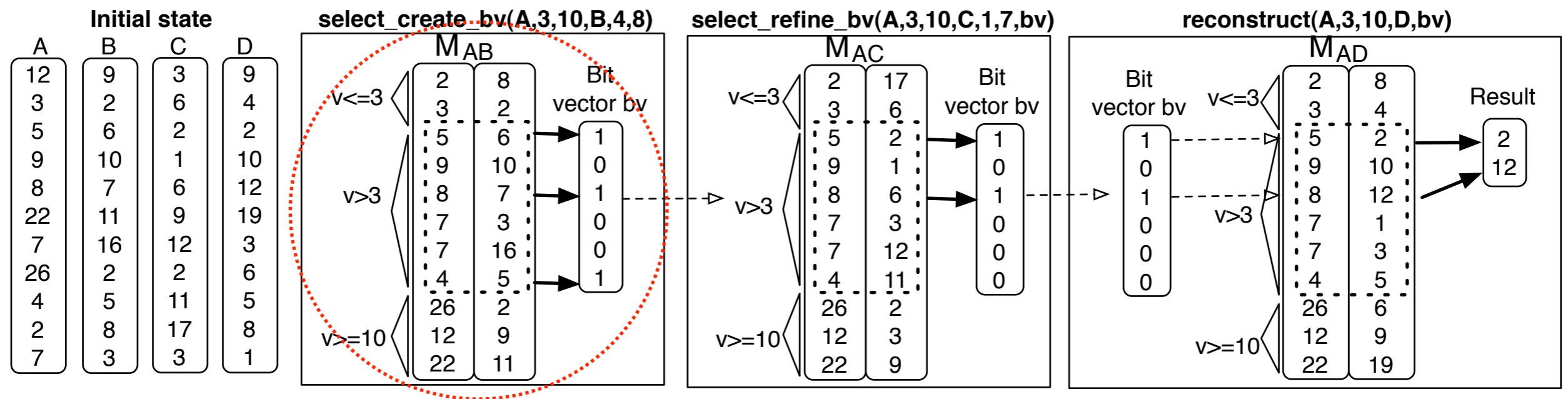
multi-selections

select D from R where $3 < A < 10$ and $4 < B < 8$ and $1 < C < 7$



multi-selections

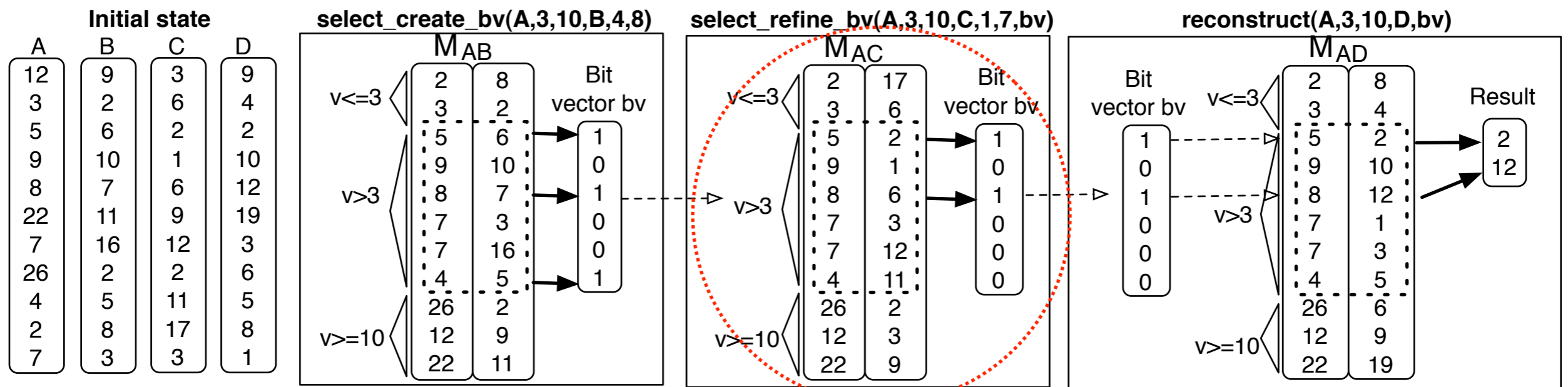
select D from R where $3 < A < 10$ and $4 < B < 8$ and $1 < C < 7$



Crack $3 < A < 10$
 Analyze tail $4 < B < 8$
 Create bit vector

multi-selections

select D from R where $3 < A < 10$ and $4 < B < 8$ and $1 < C < 7$

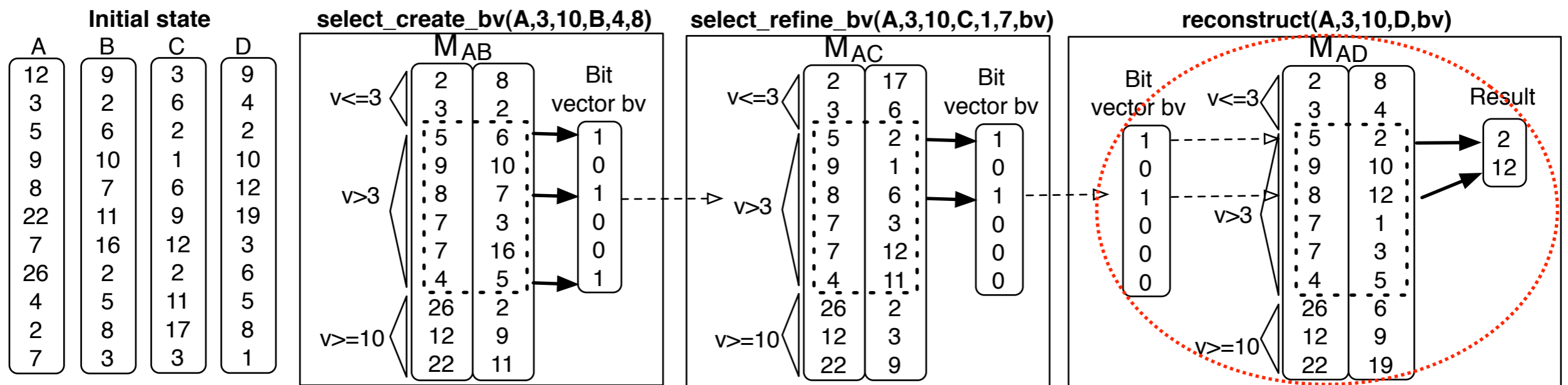


Crack $3 < A < 10$
 Analyze tail $4 < B < 8$
 Create bit vector

Align $3 < A < 10$
 Analyze tail $1 < C < 7$
 Refine bit vector

multi-selections

select D from R where $3 < A < 10$ and $4 < B < 8$ and $1 < C < 7$



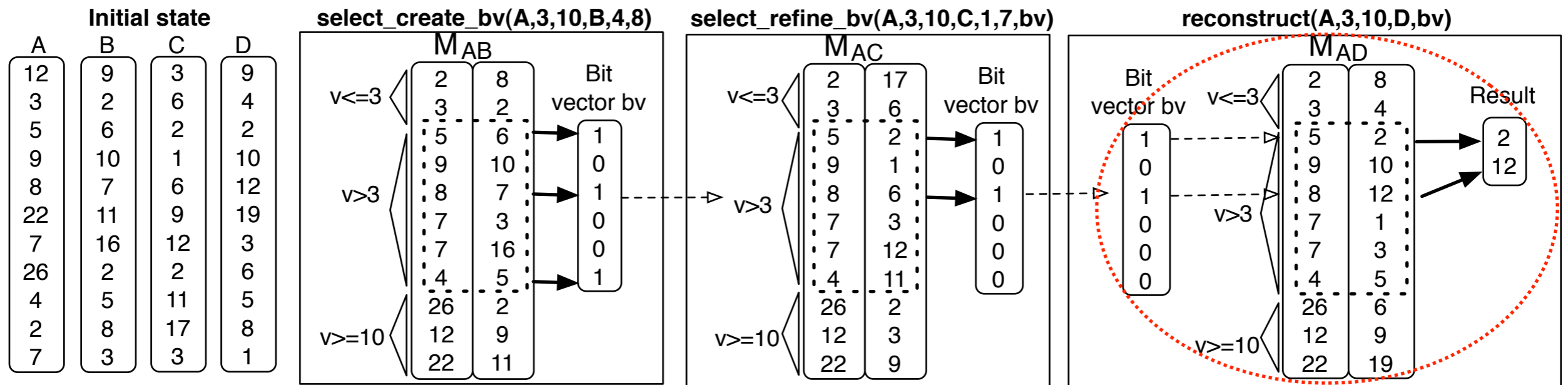
Crack $3 < A < 10$
Analyze tail $4 < B < 8$
Create bit vector

Align $3 < A < 10$
Analyze tail $1 < C < 7$
Refine bit vector

Align $3 < A < 10$
Grab tail values

multi-selections

select D from R where $3 < A < 10$ and $4 < B < 8$ and $1 < C < 7$

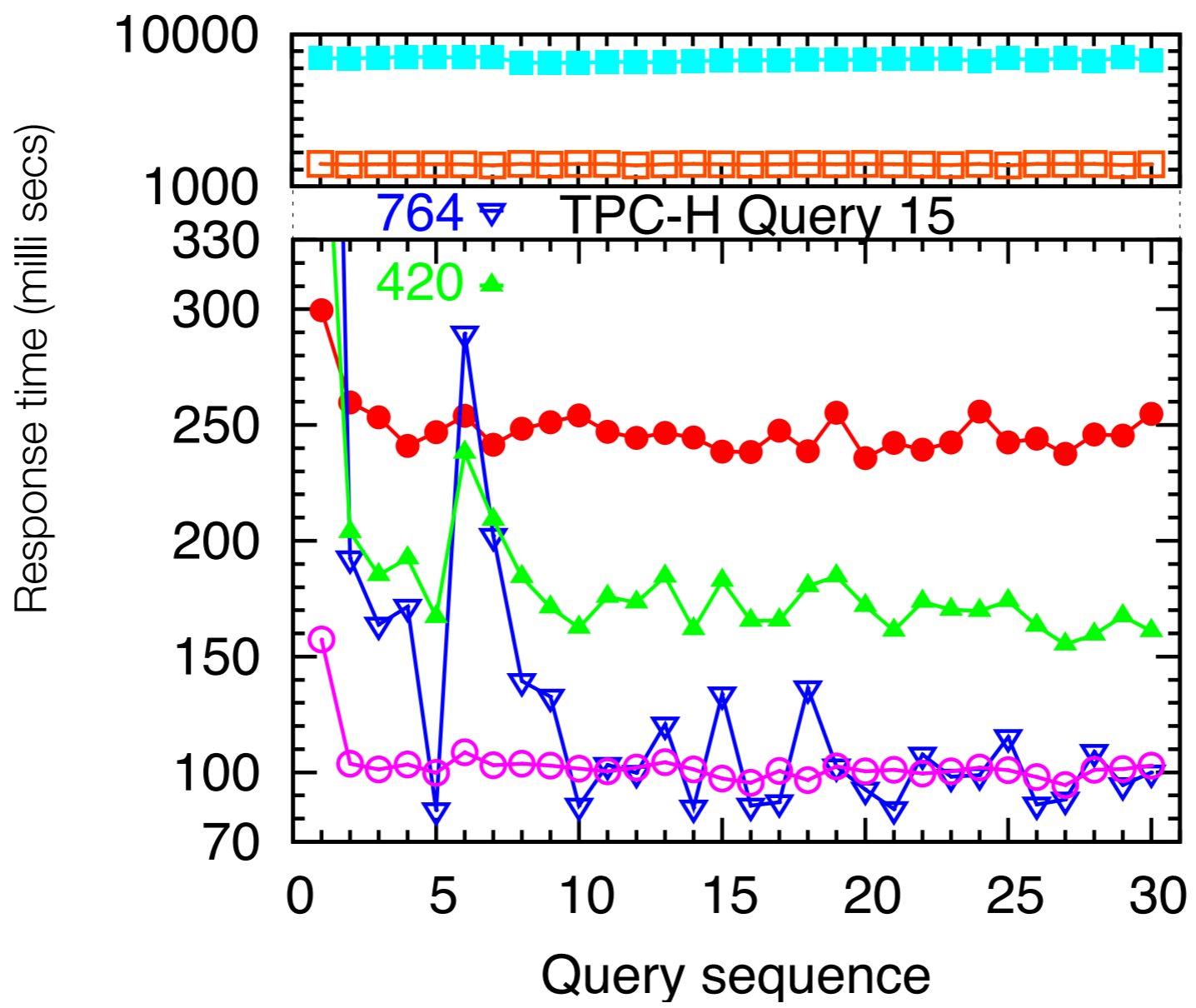
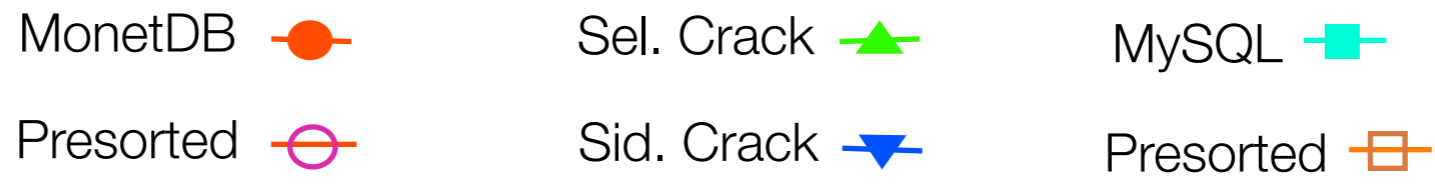


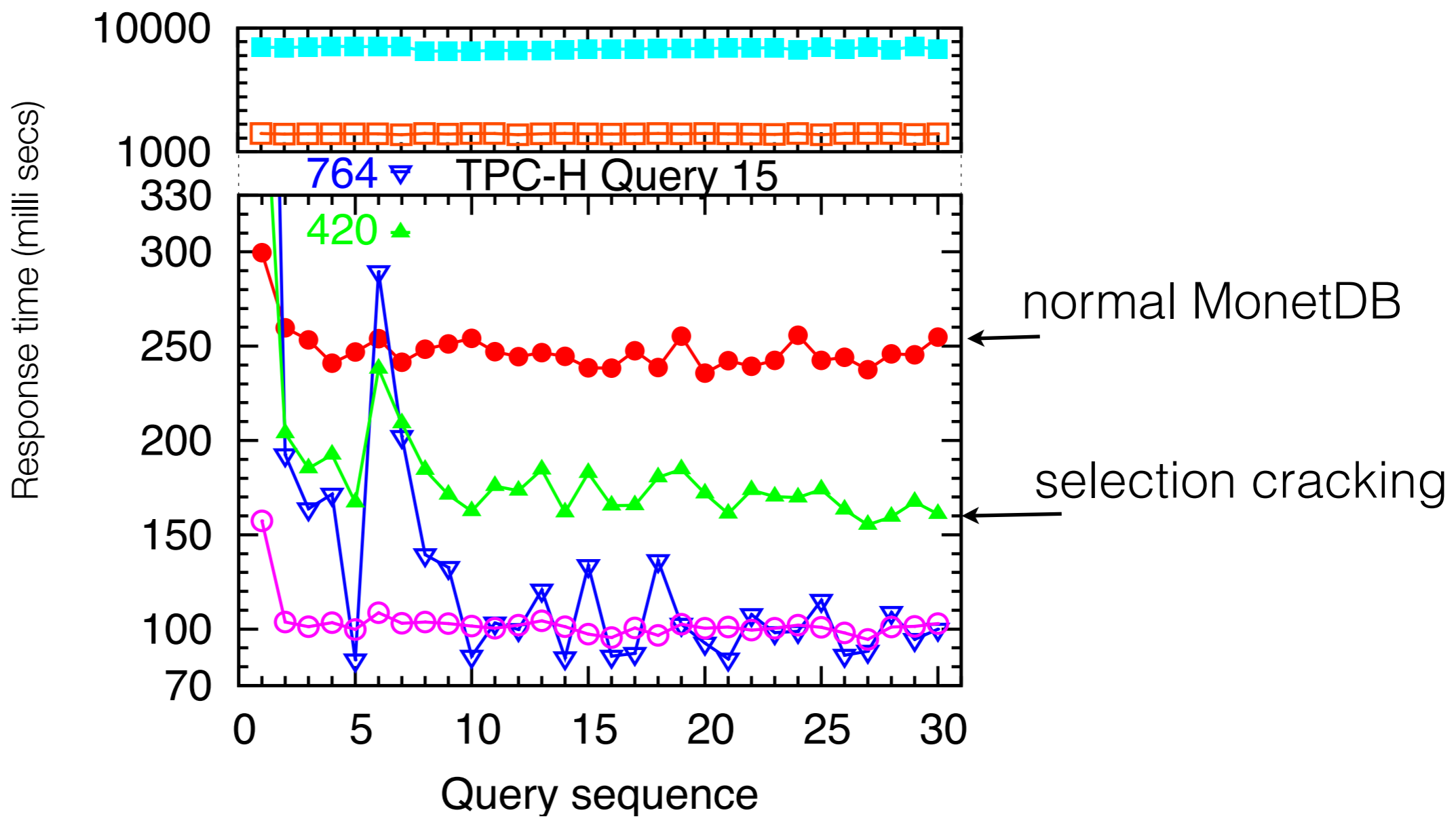
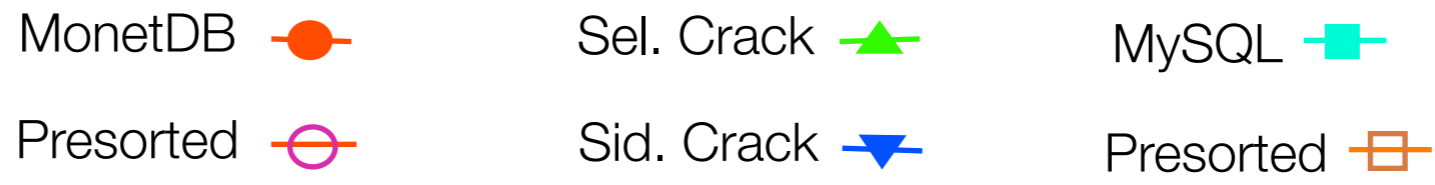
Crack $3 < A < 10$
 Analyze tail $4 < B < 8$
 Create bit vector

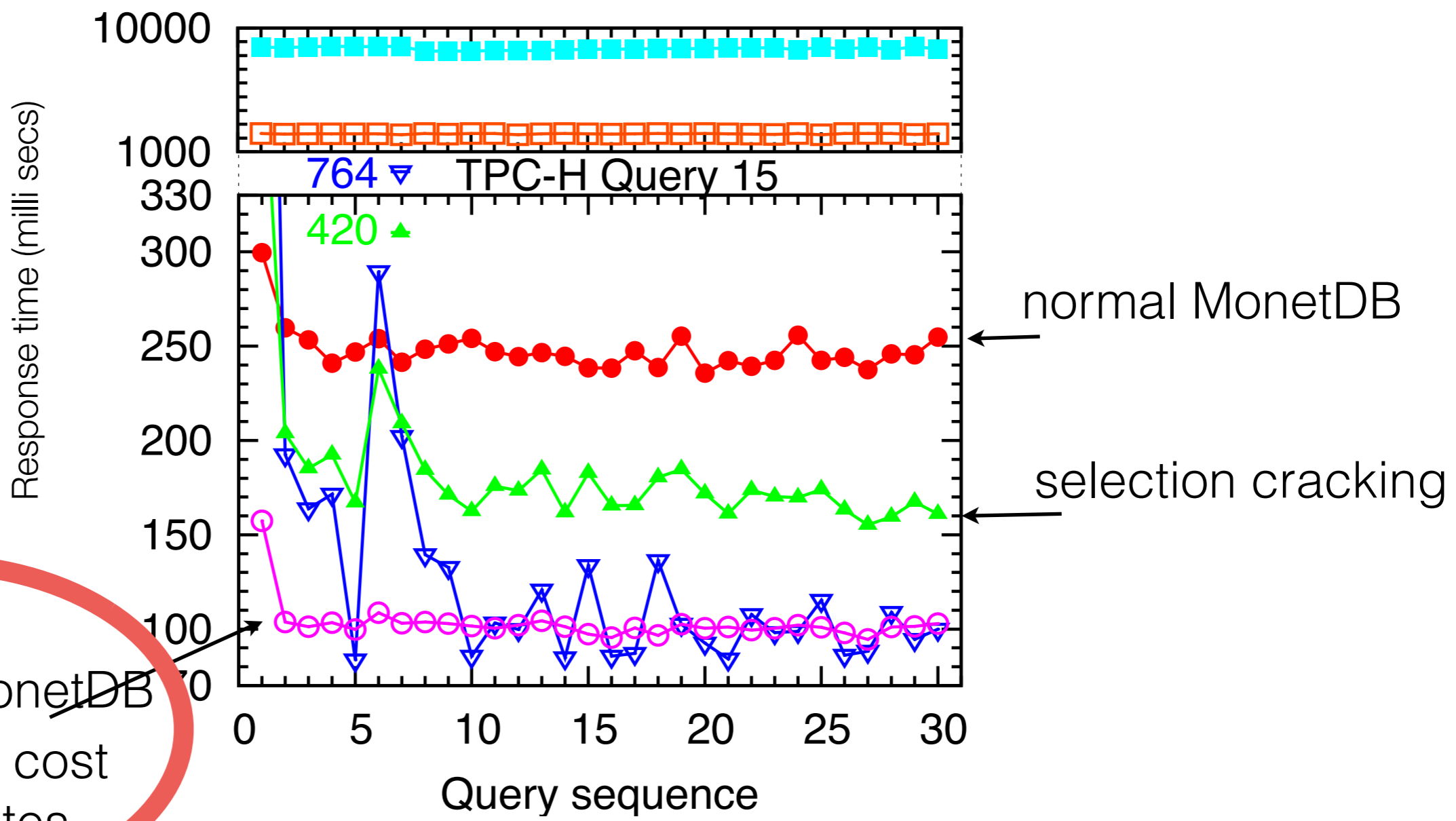
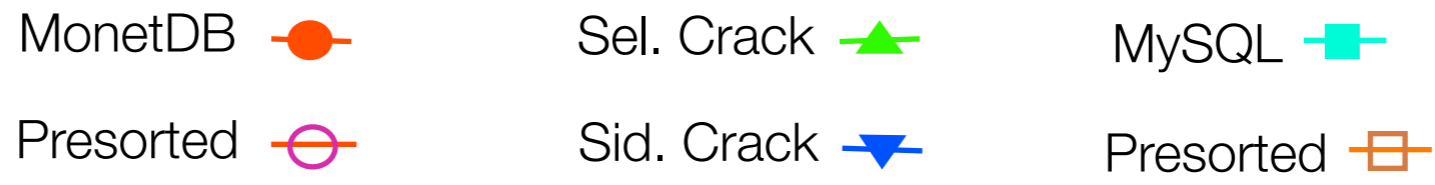
Align $3 < A < 10$
 Analyze tail $1 < C < 7$
 Refine bit vector

Align $3 < A < 10$
 Grab tail values

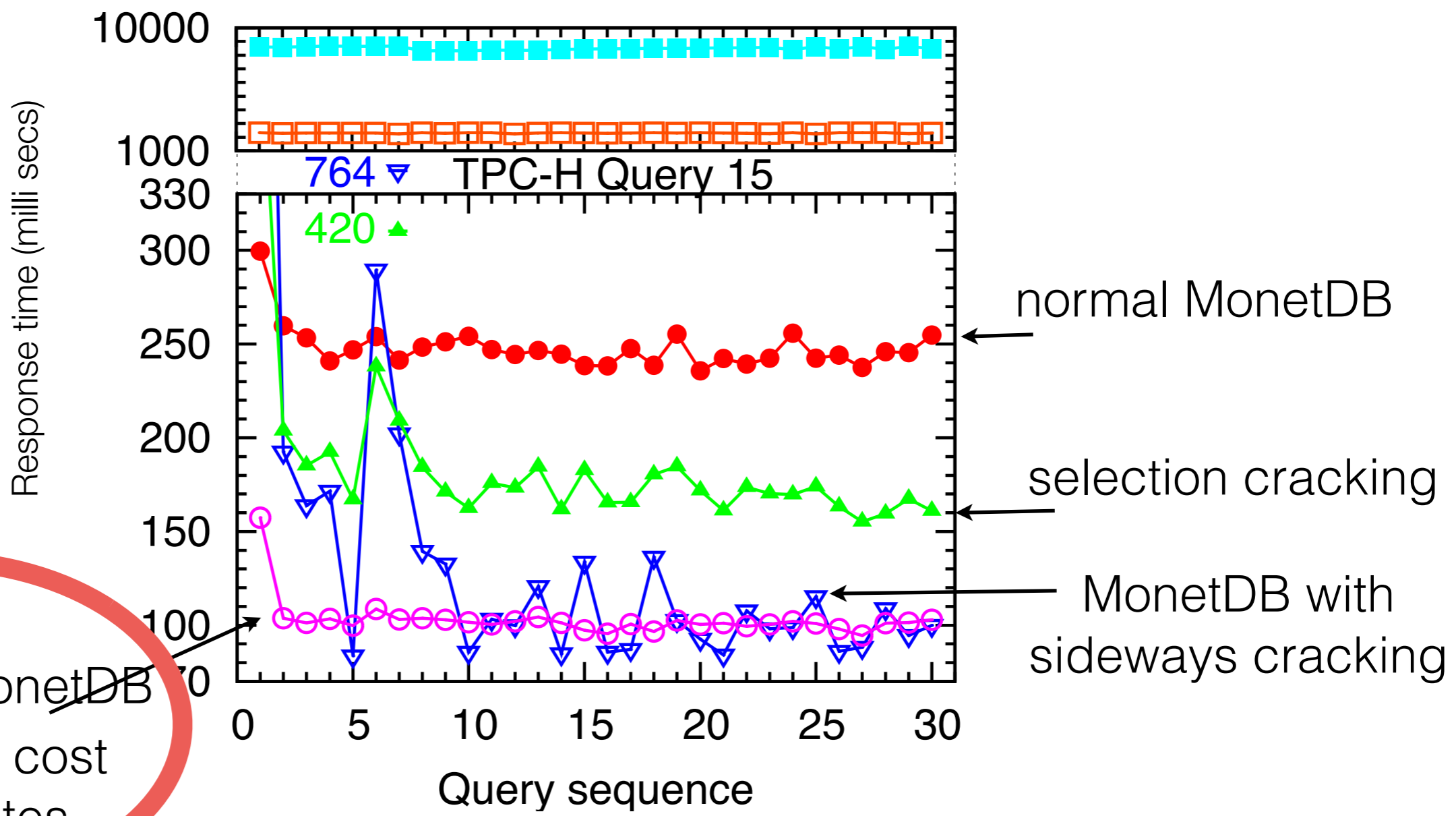
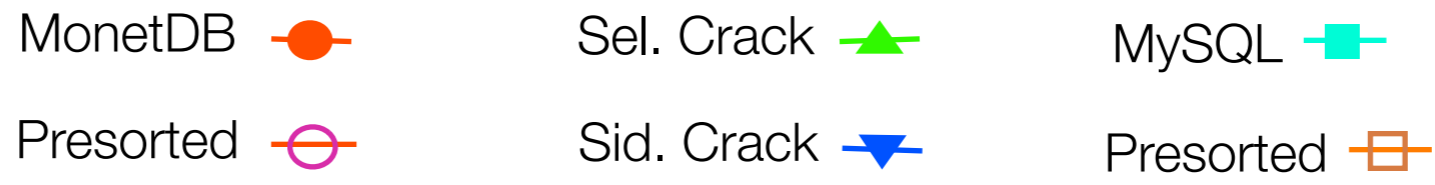
Use histogram-like info from maps to choose map set



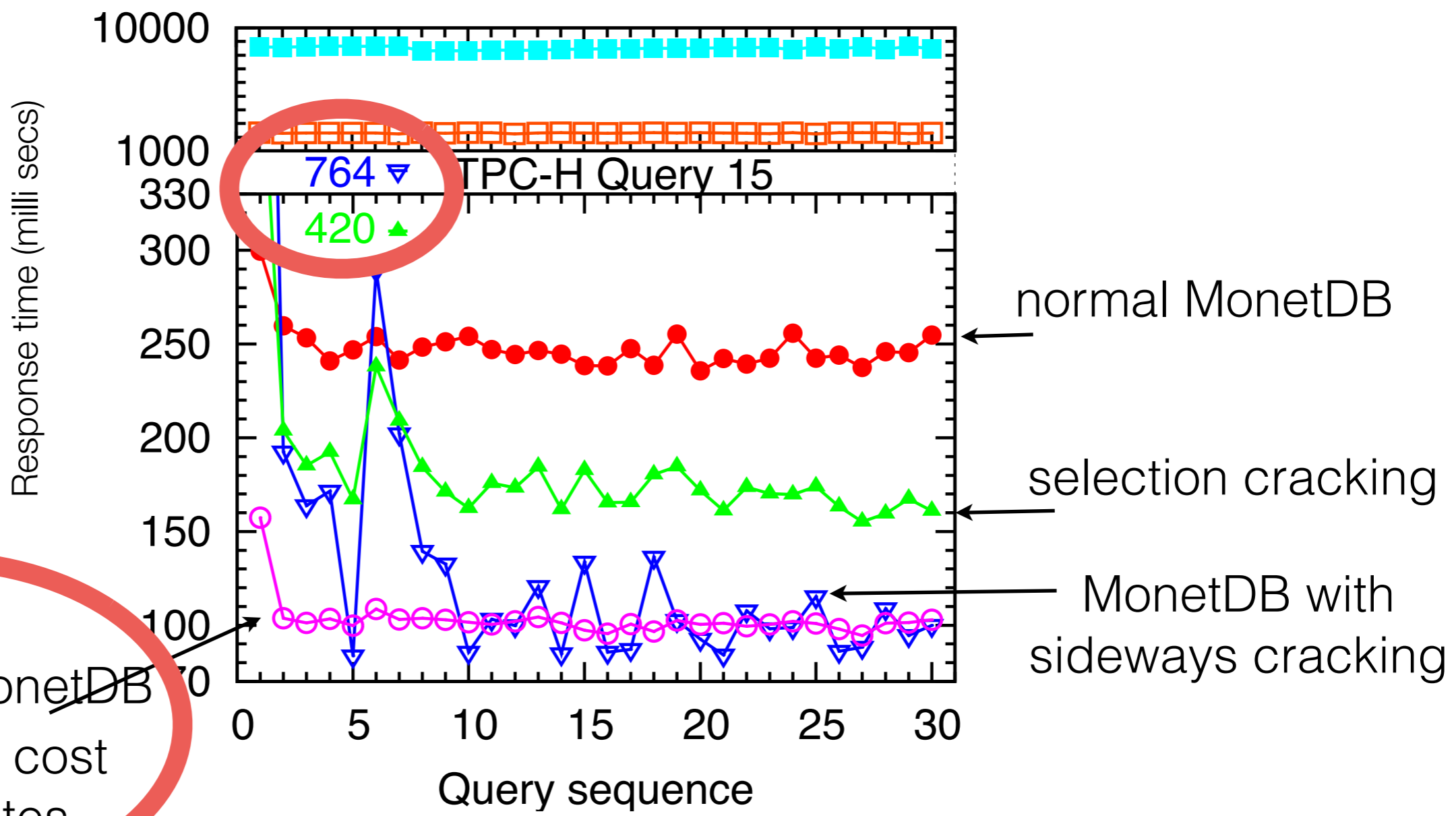
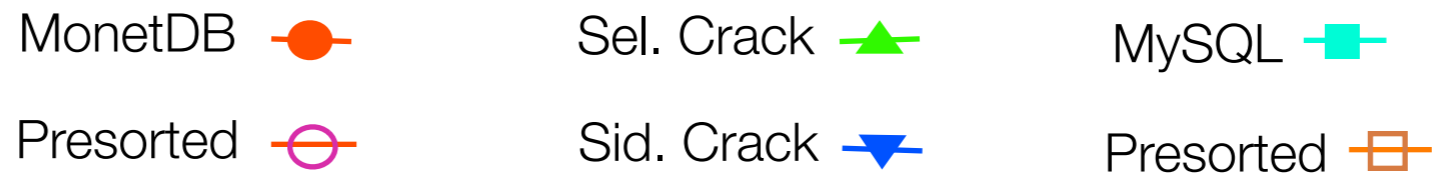




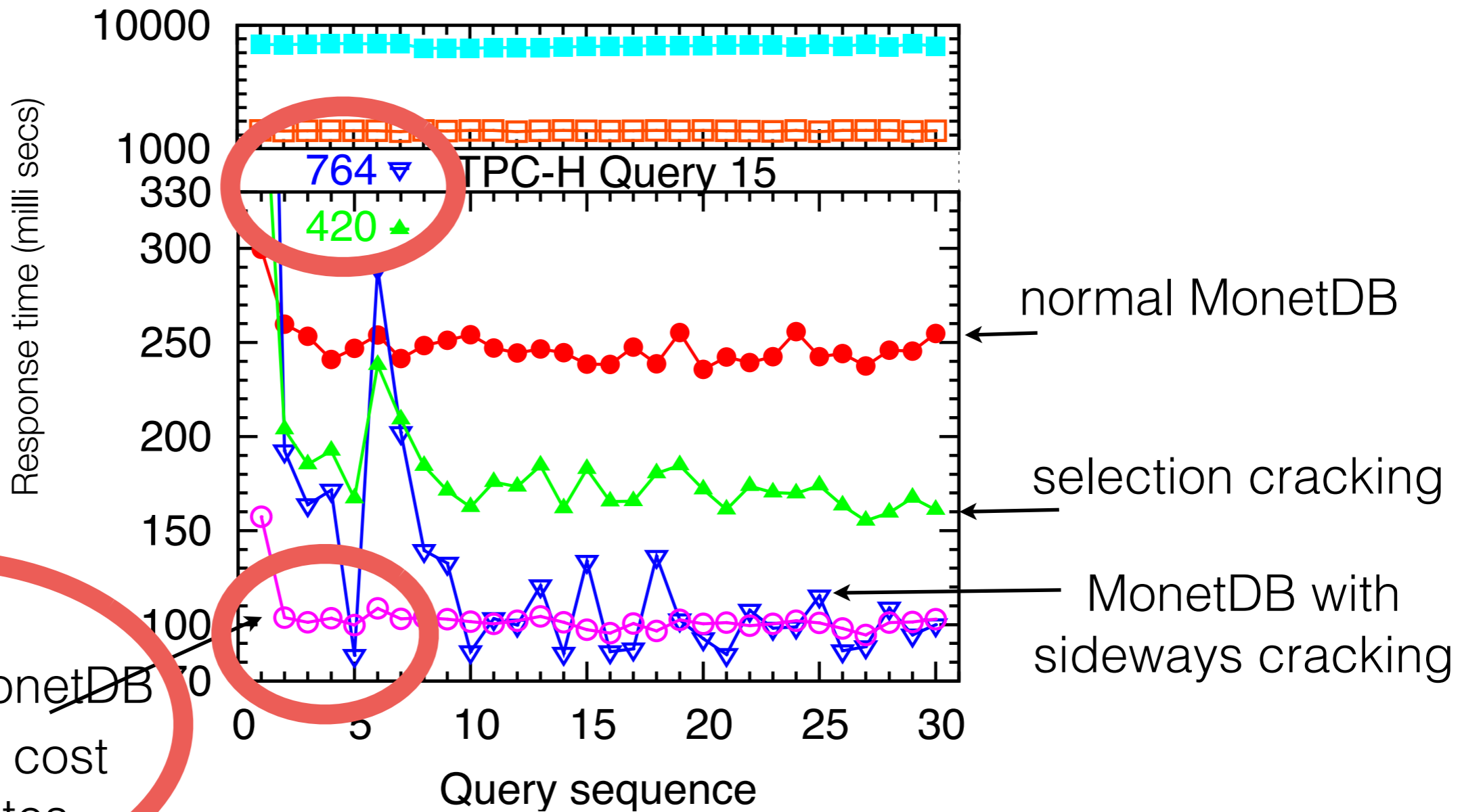
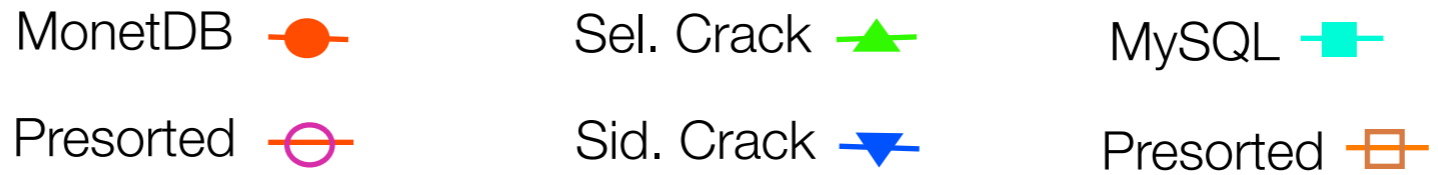
presorted MonetDB
preparation cost
3-14 minutes



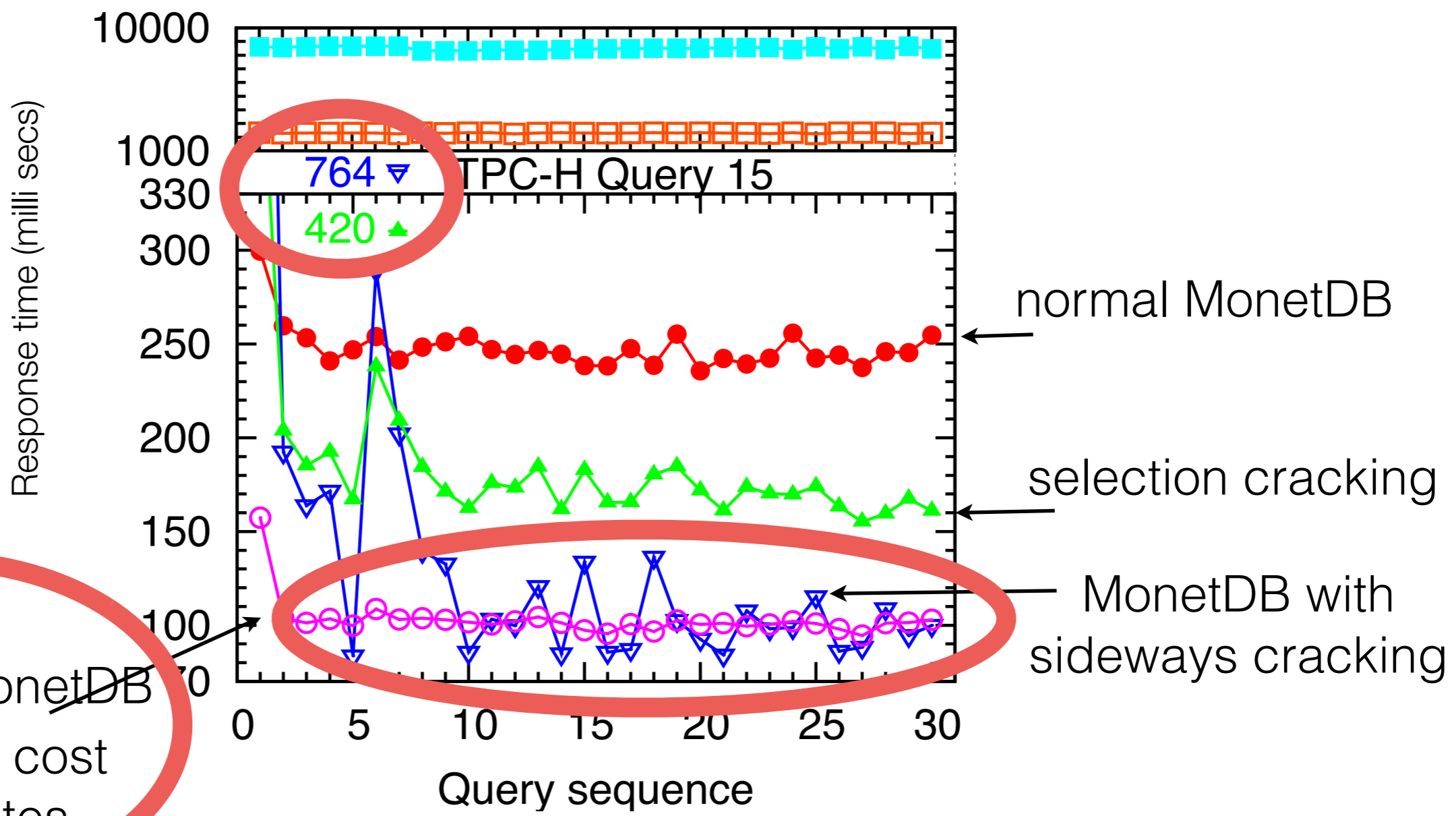
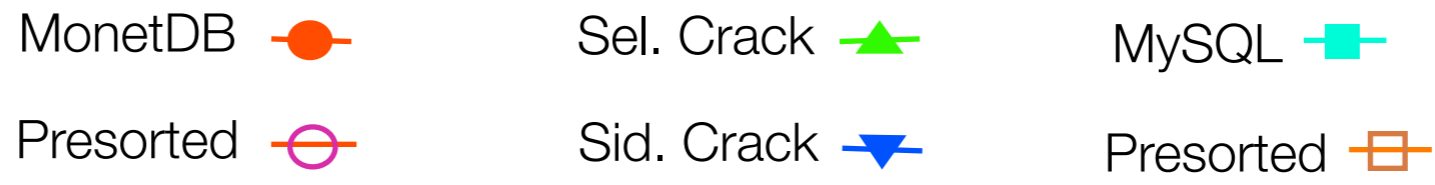
presorted MonetDB
preparation cost
3-14 minutes



presorted MonetDB
preparation cost
3-14 minutes



presorted MonetDB
preparation cost
3-14 minutes



presorted MonetDB
preparation cost
3-14 minutes

updates

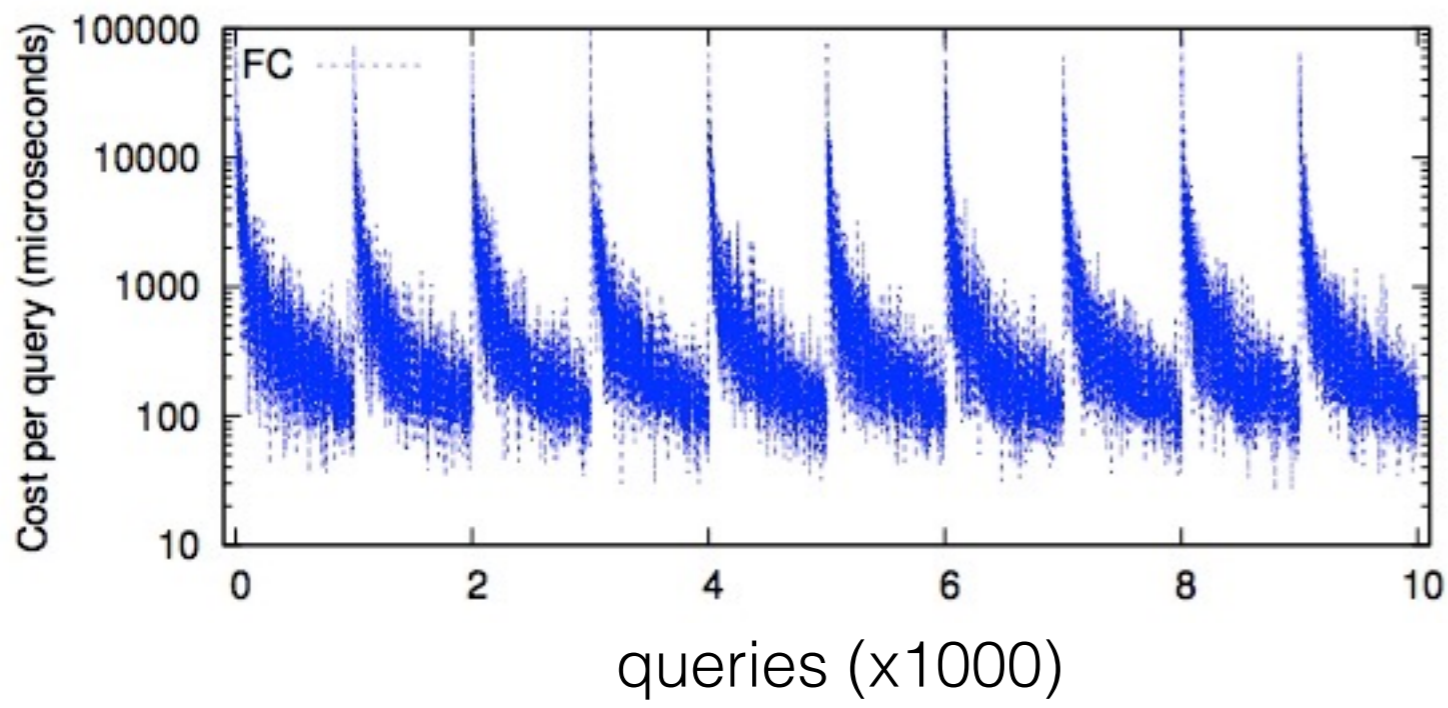
cracking indices are auxiliary data structures
can be dropped any time

column of 10M tuples, random queries
1000 random insertions every 1000 queries

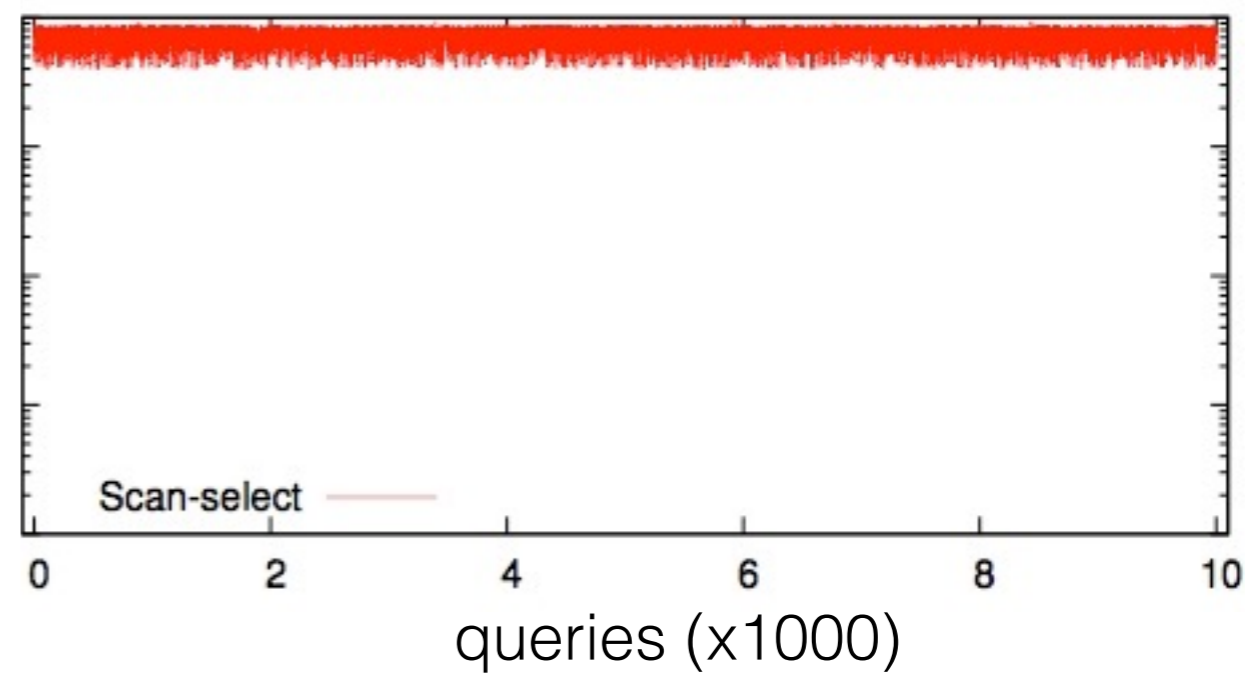
forgetting

when updates arrive, drop the index

crack (forget)



scan



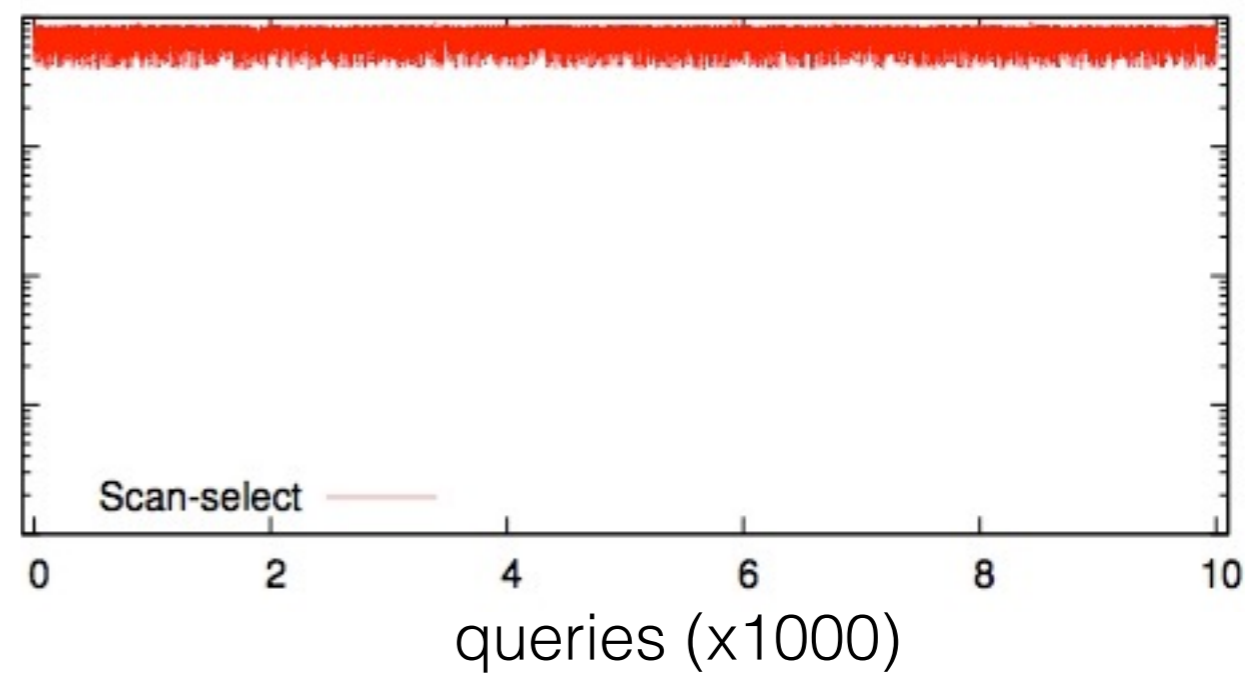
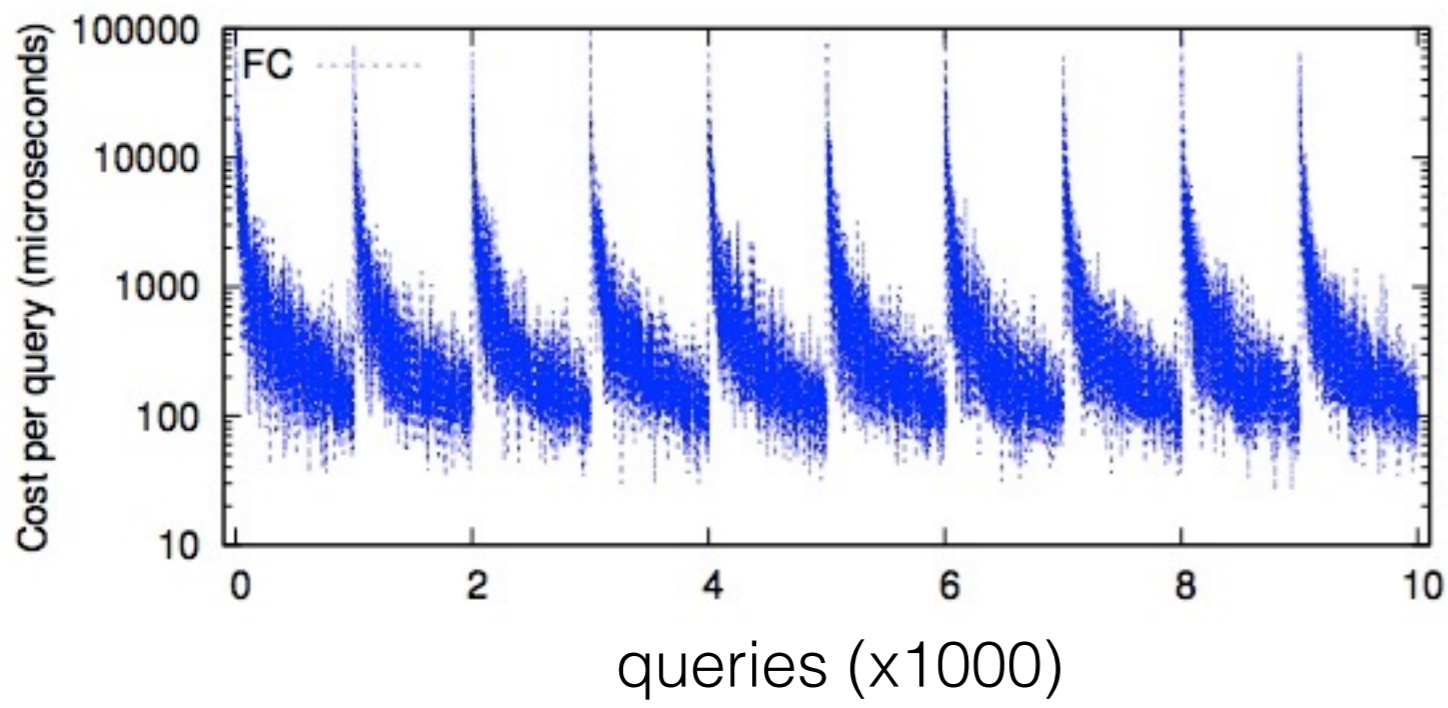
column of 10M tuples, random queries
1000 random insertions every 1000 queries

forgetting

when updates arrive, drop the index

crack (forget)

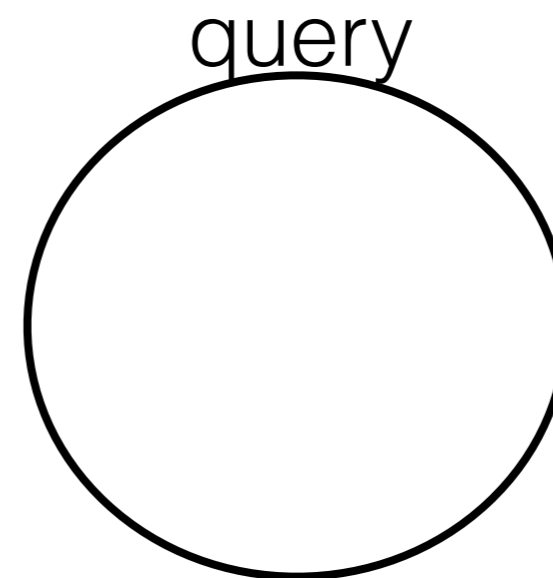
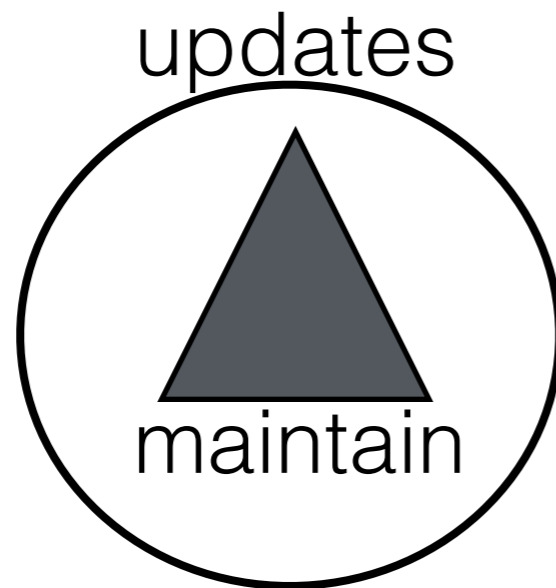
scan



we do not exploit past cracking

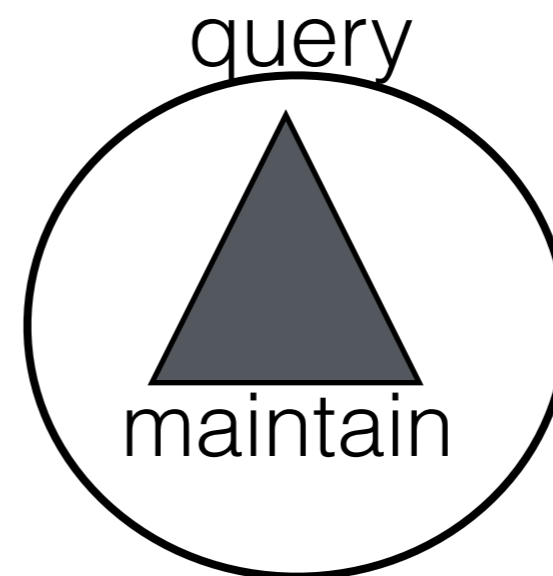
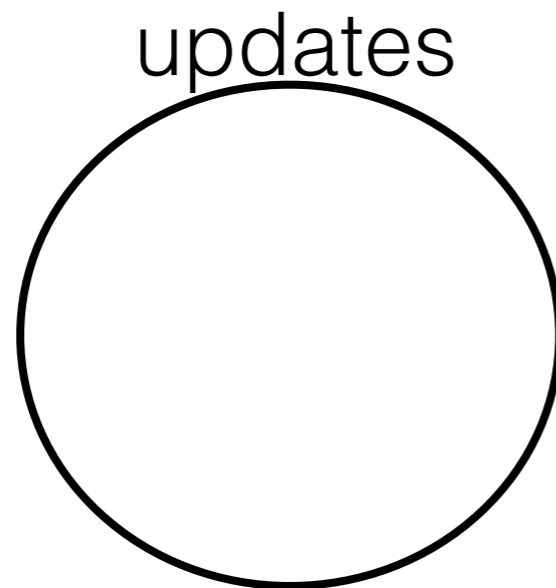
cracking updates

log updates and apply on-line and on-demand during cracking



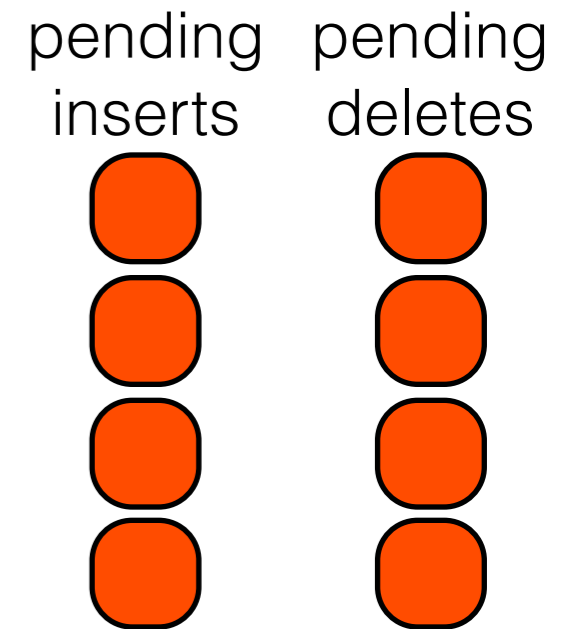
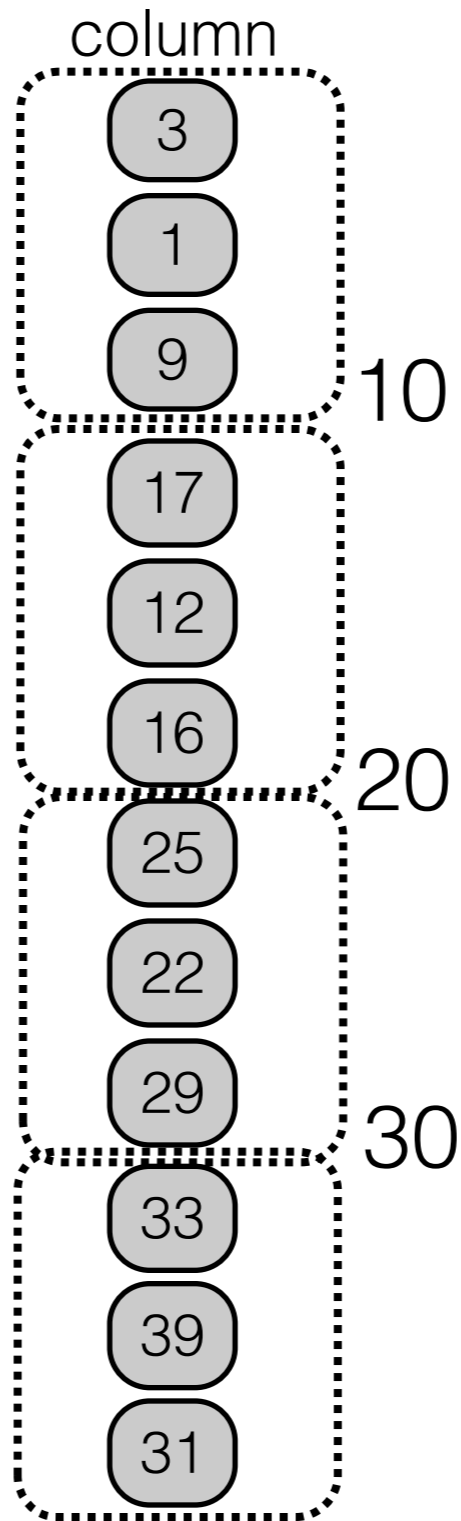
cracking updates

log updates and apply on-line and on-demand during cracking



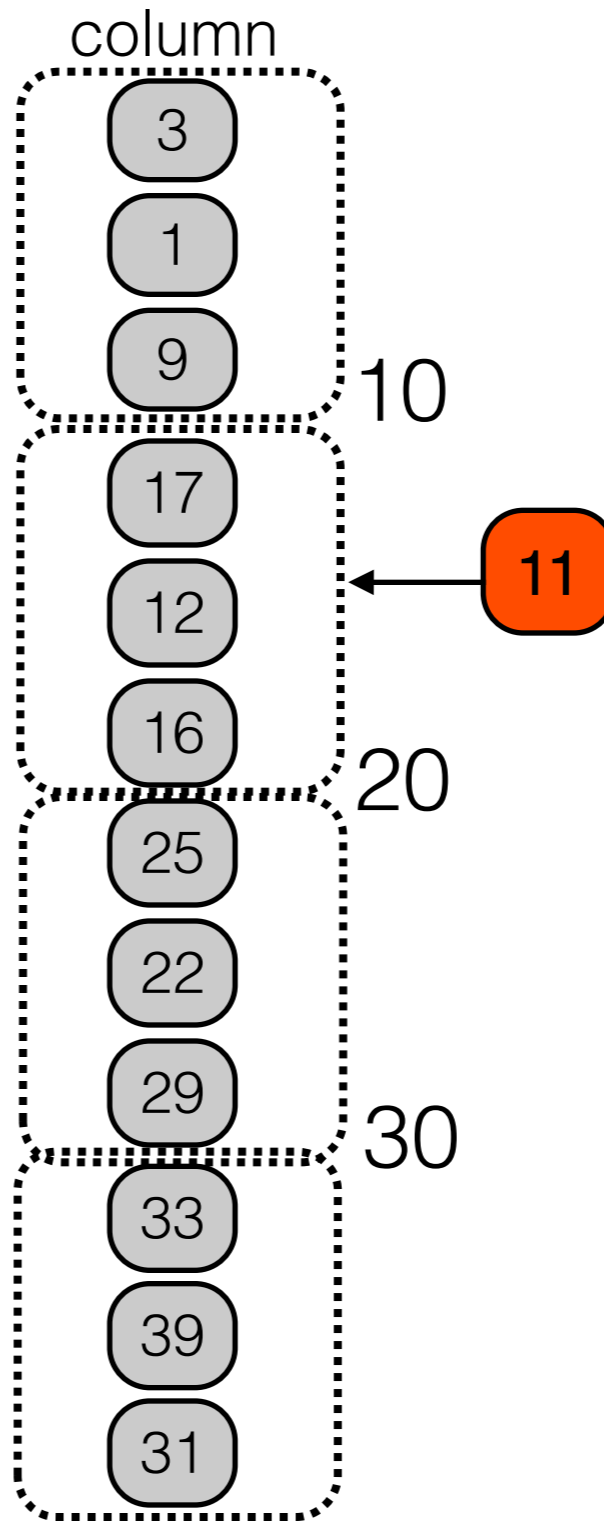
goal: minimize physical actions

- arrays are dense
- pieces are ordered
- values in a piece are not ordered

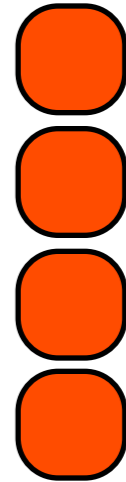


goal: minimize physical actions

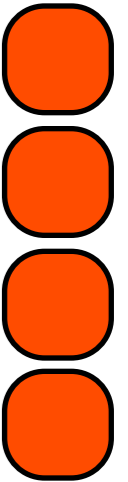
- arrays are dense
- pieces are ordered
- values in a piece are not ordered



pending
inserts

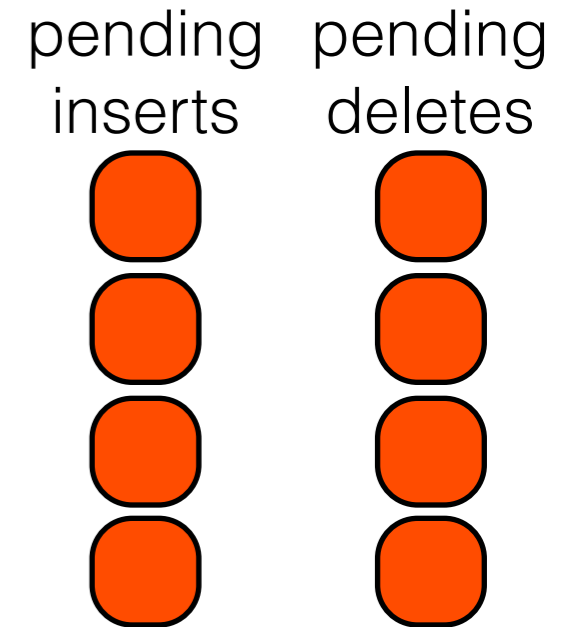
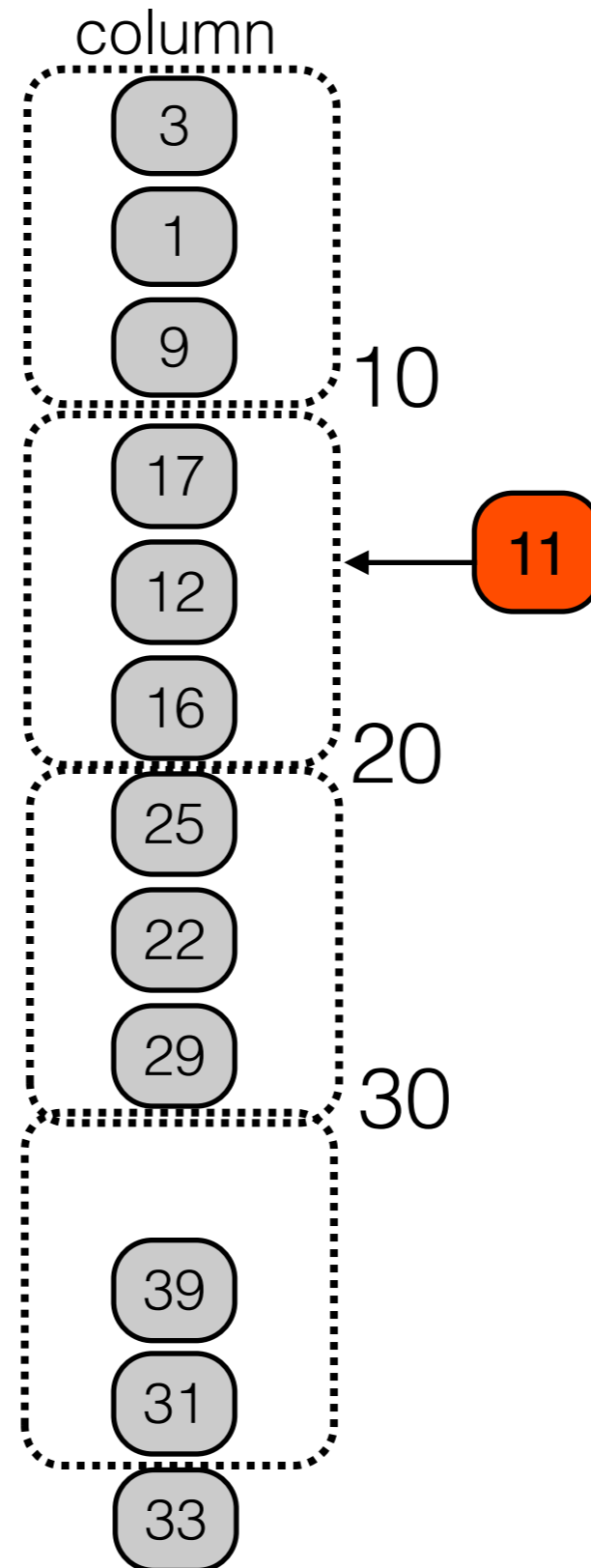


pending
deletes



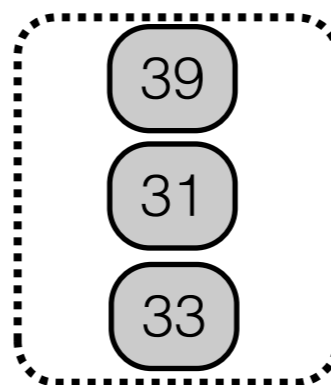
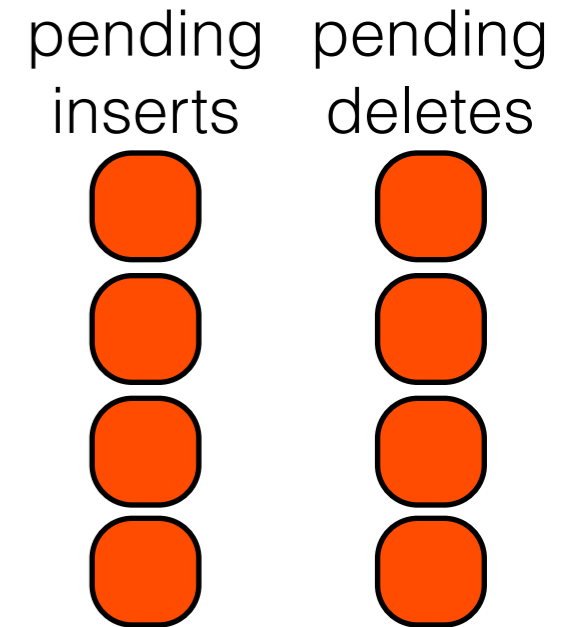
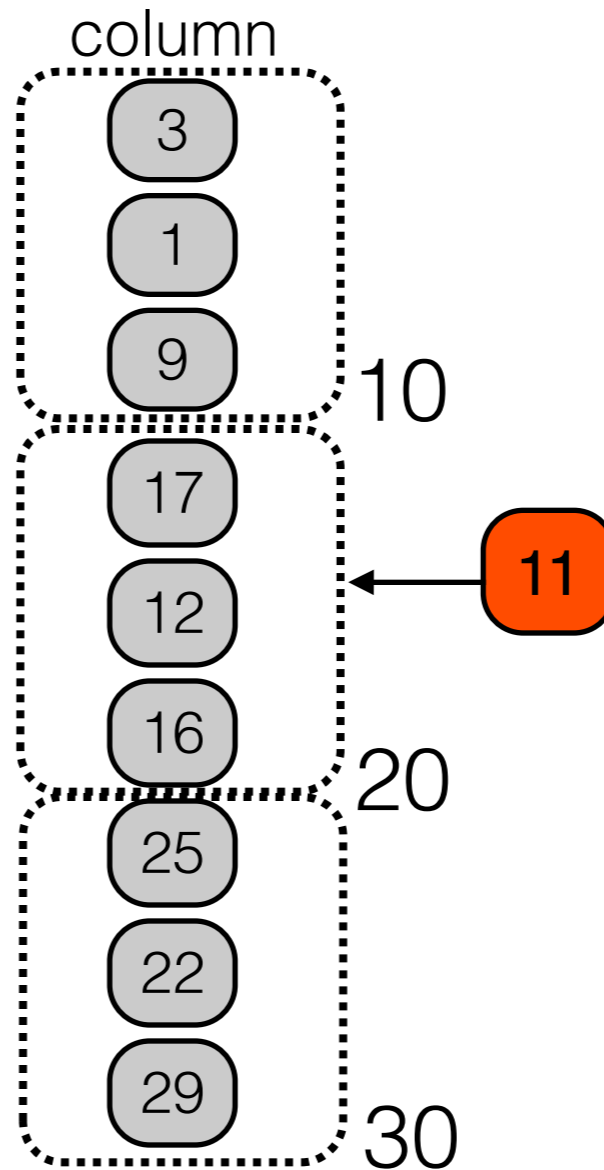
goal: minimize physical actions

- arrays are dense
- pieces are ordered
- values in a piece are not ordered



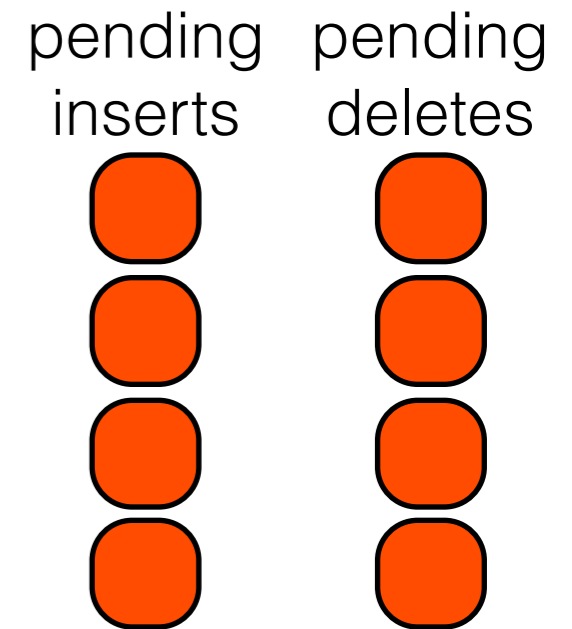
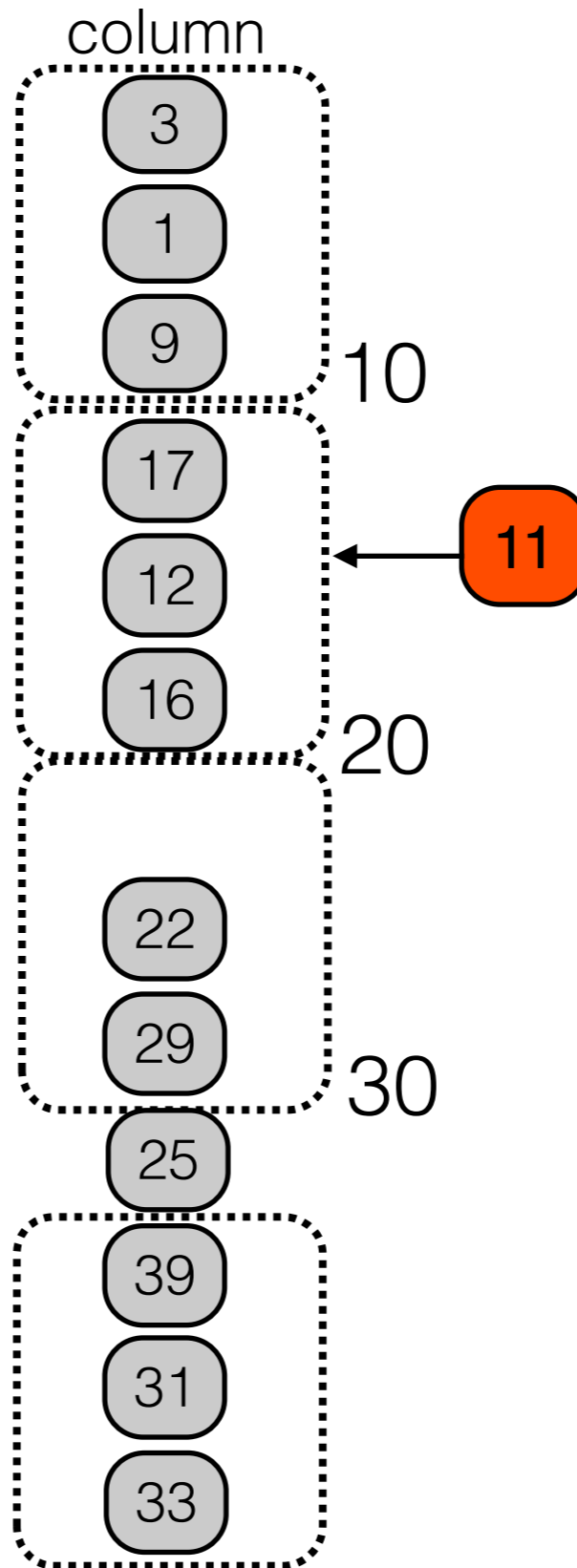
goal: minimize physical actions

- arrays are dense
- pieces are ordered
- values in a piece are not ordered



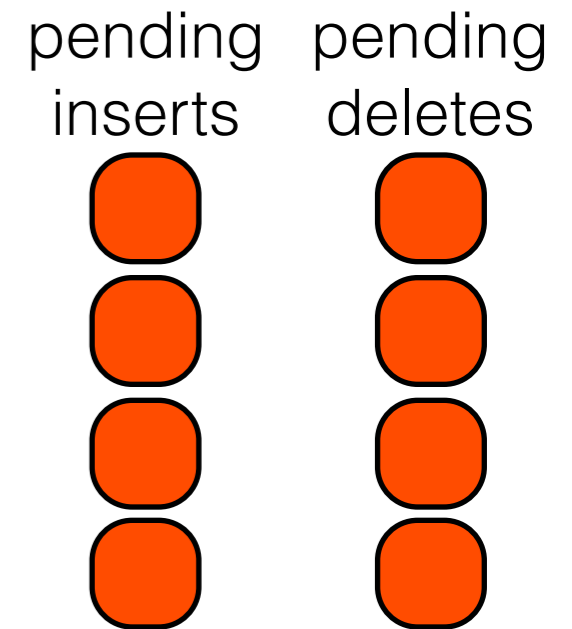
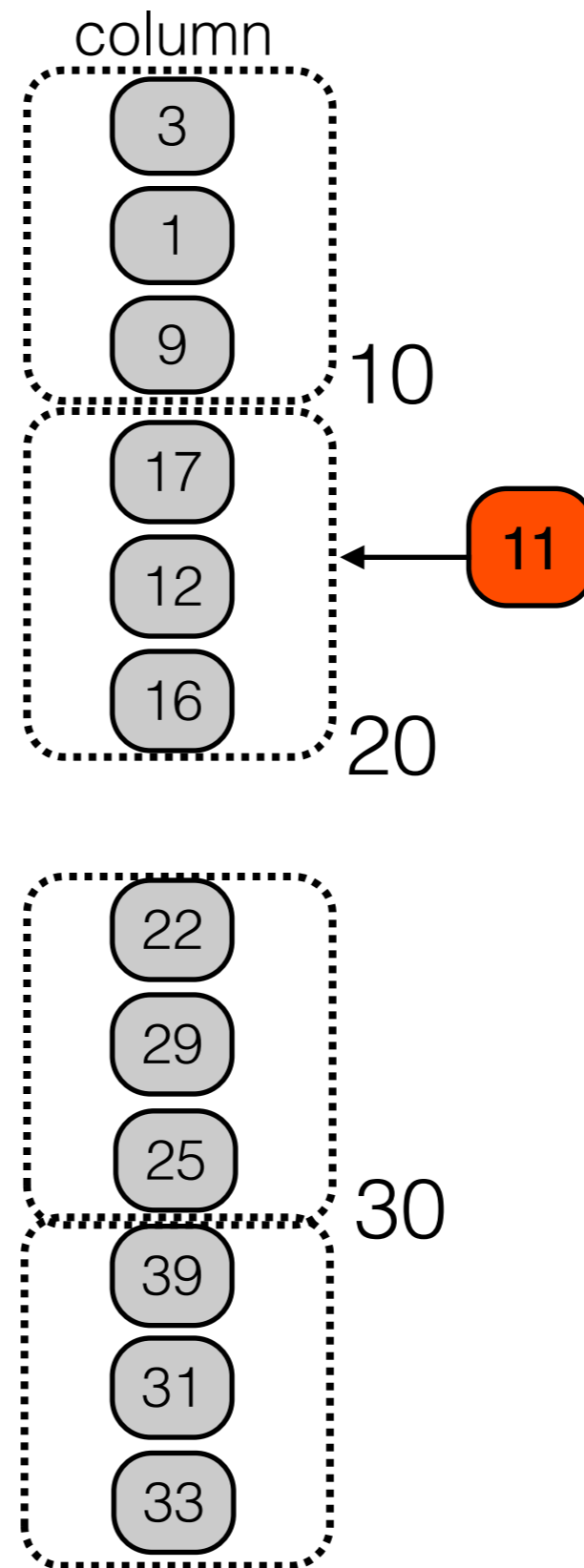
goal: minimize physical actions

- arrays are dense
- pieces are ordered
- values in a piece are not ordered



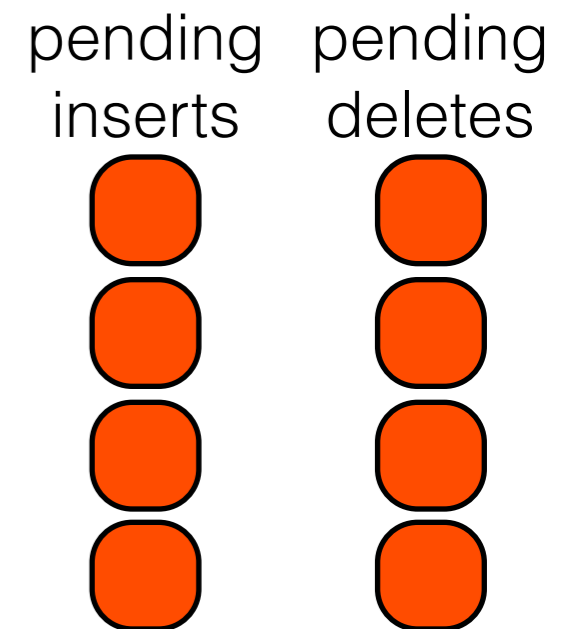
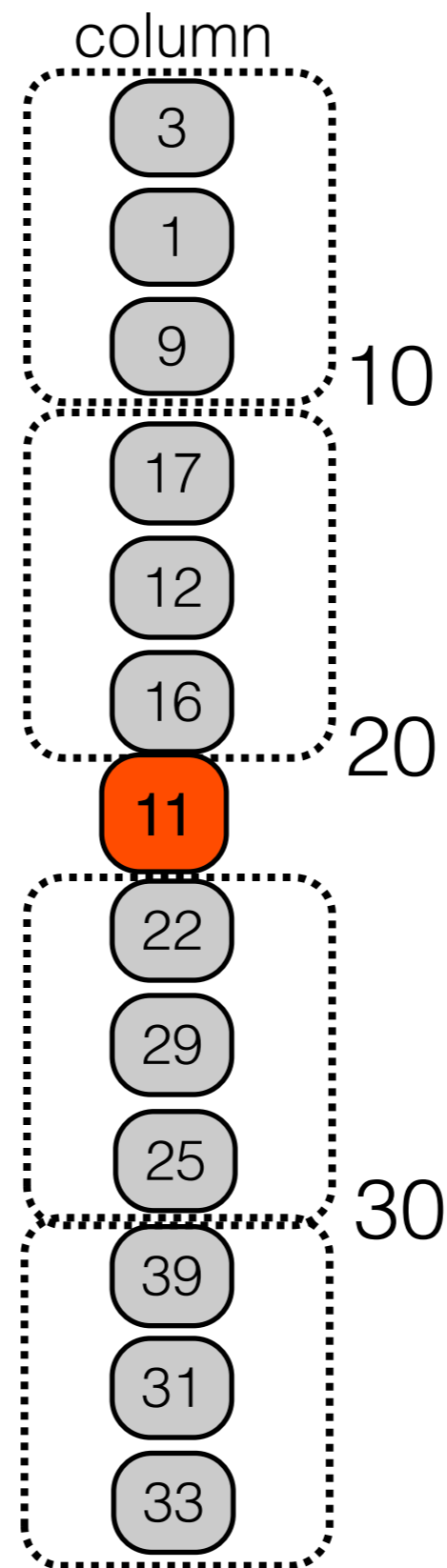
goal: minimize physical actions

- arrays are dense
- pieces are ordered
- values in a piece are not ordered



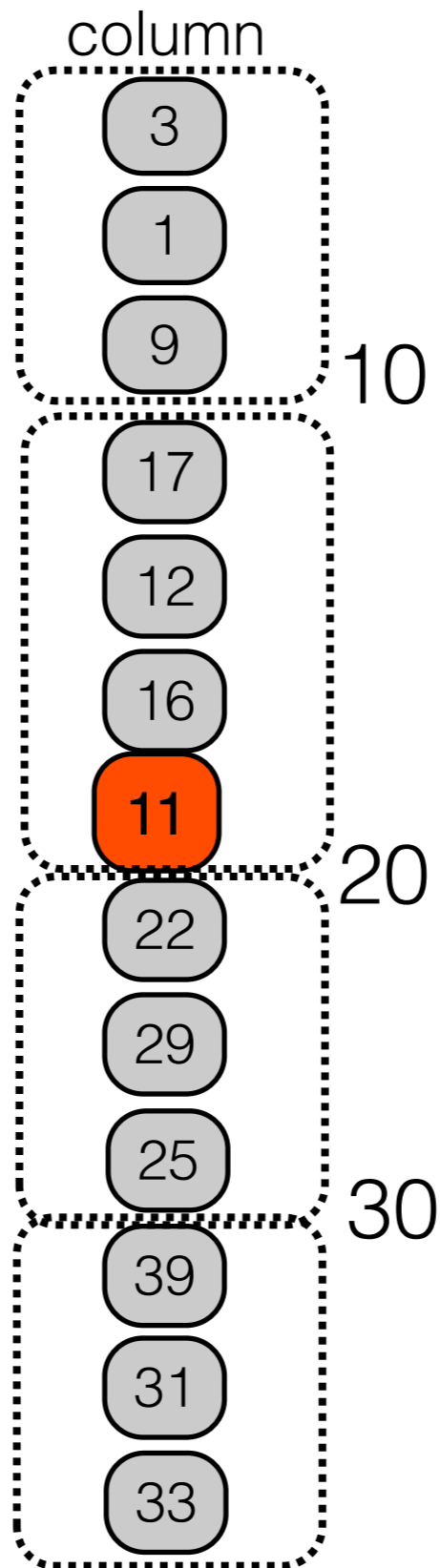
goal: minimize physical actions

- arrays are dense
- pieces are ordered
- values in a piece are not ordered

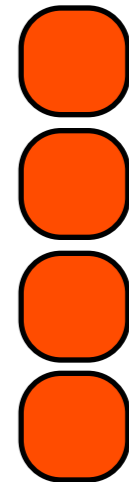


goal: minimize physical actions

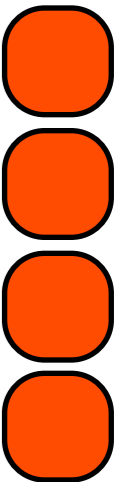
- arrays are dense
- pieces are ordered
- values in a piece are not ordered



pending
inserts

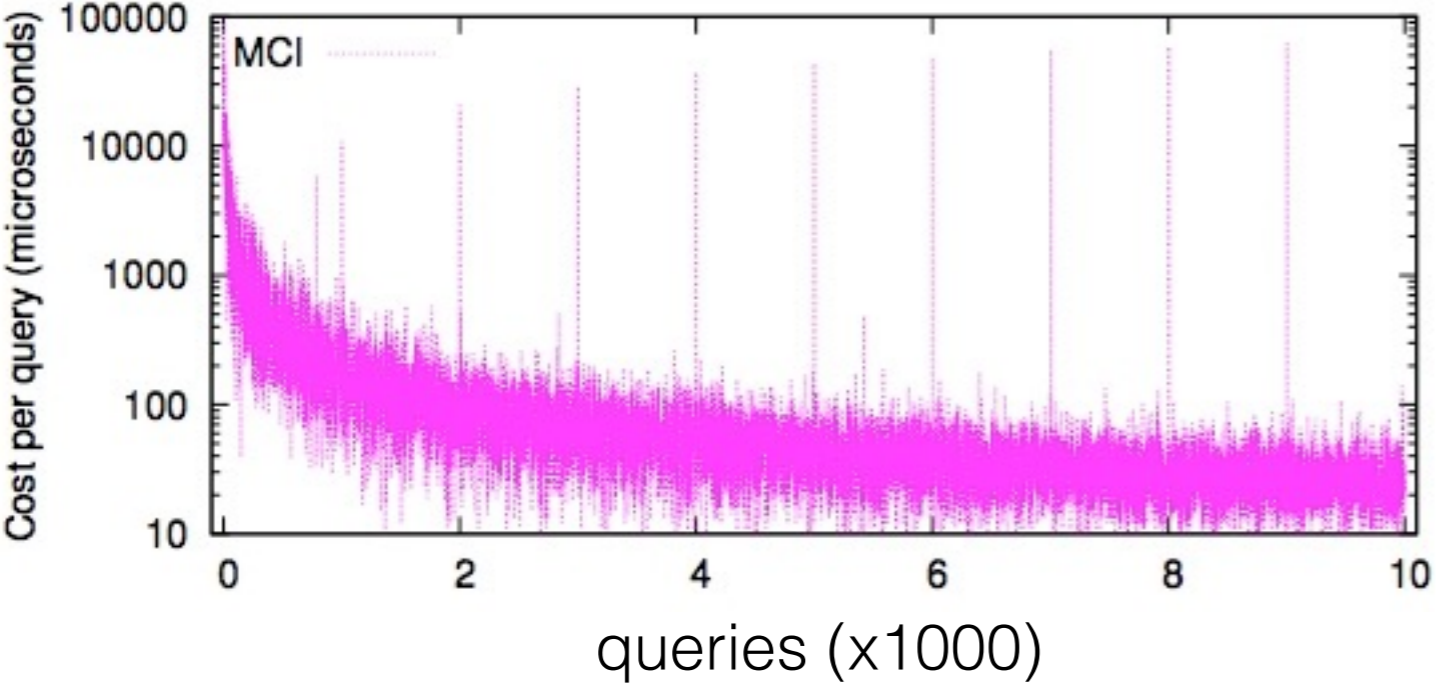


pending
deletes

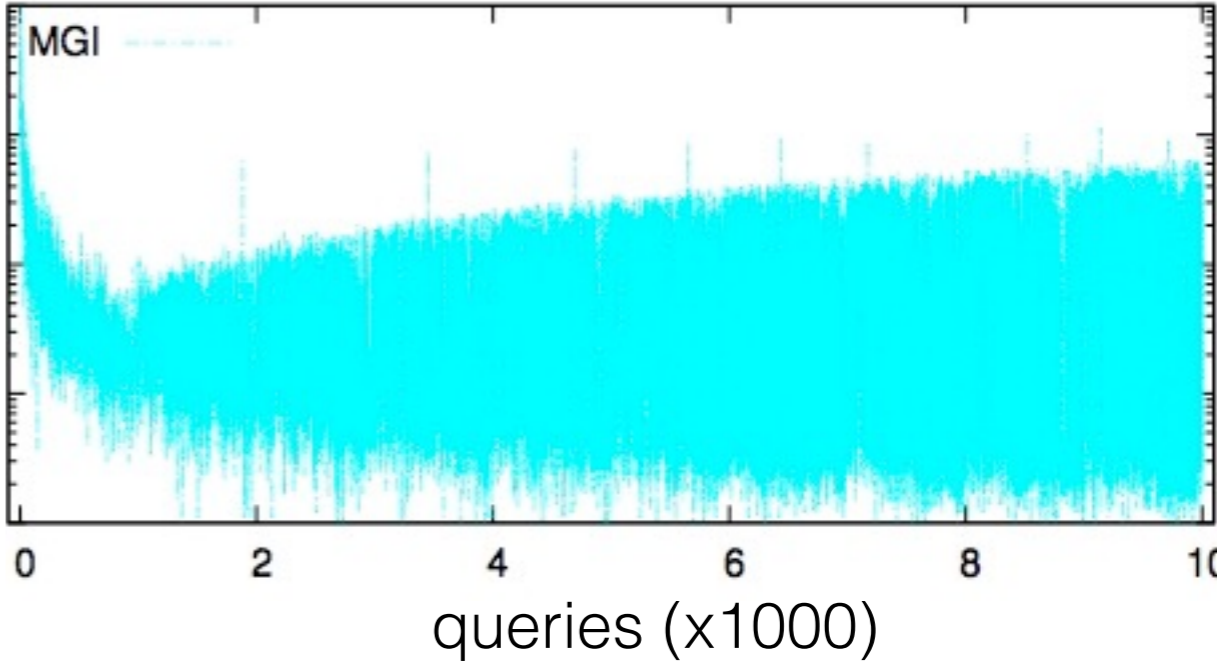


column of 10M tuples, random queries
1000 random insertions every 1000 queries

merge all updates

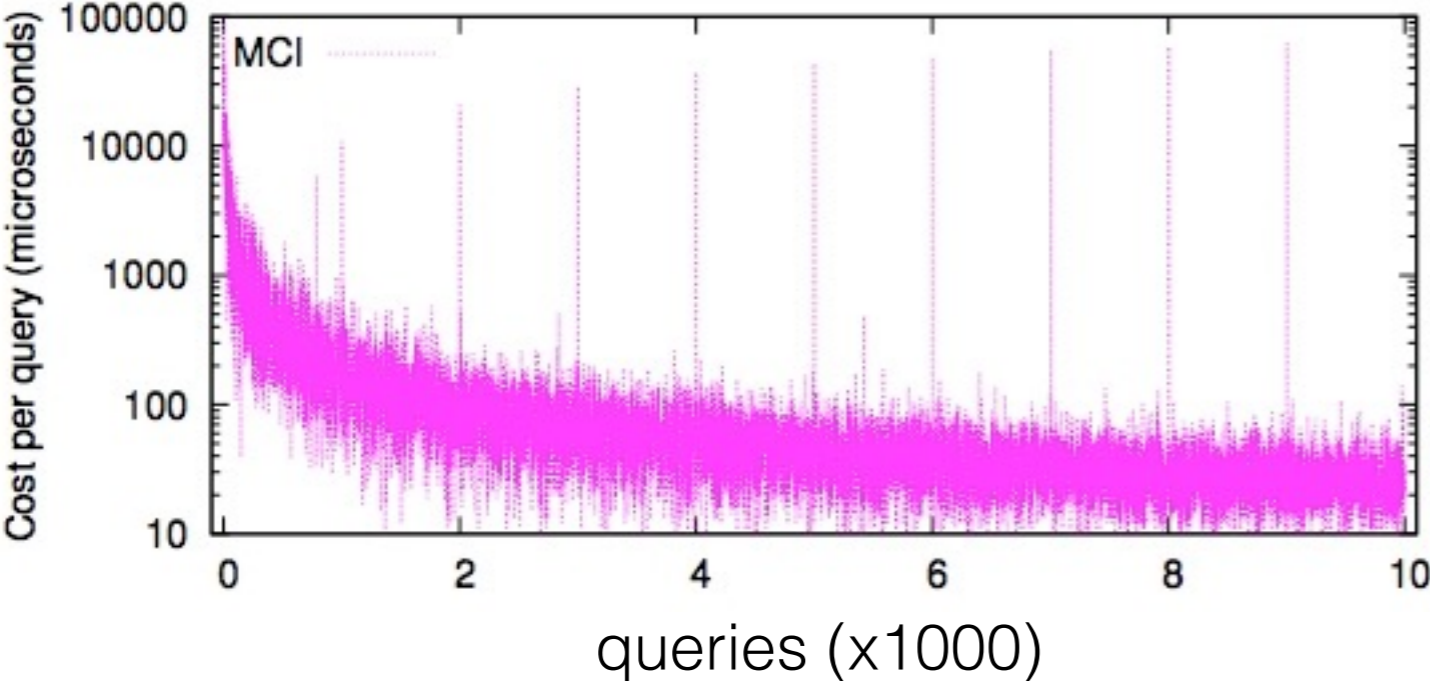


merge only qualifying updates

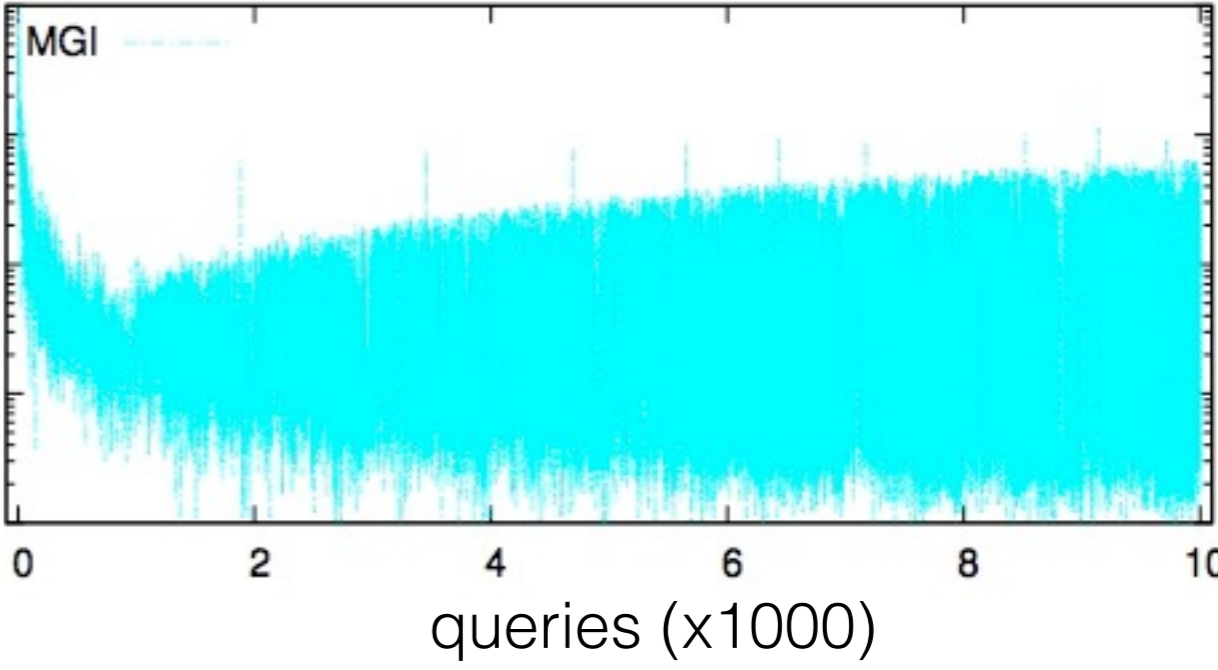


column of 10M tuples, random queries
1000 random insertions every 1000 queries

merge all updates

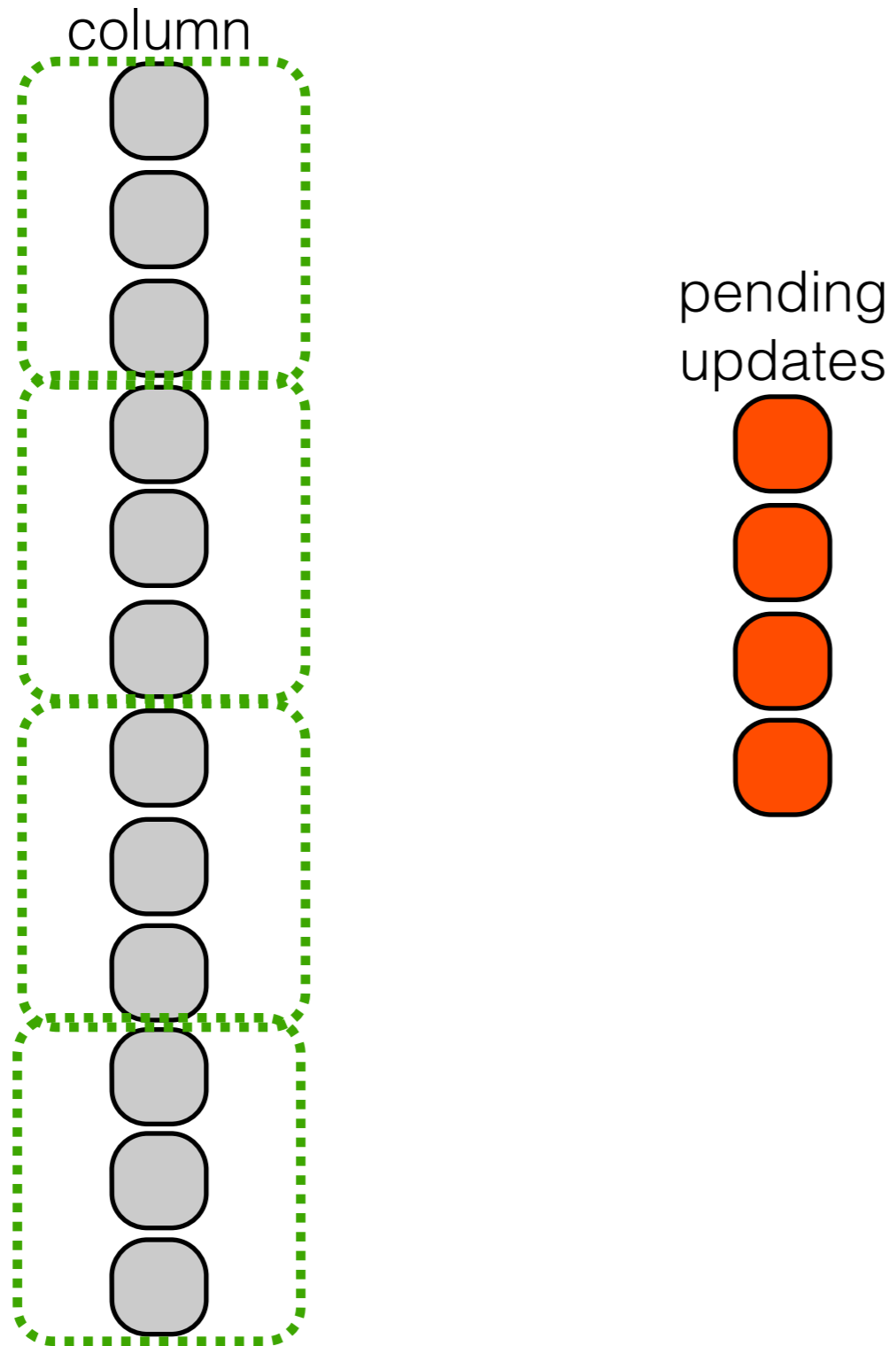


merge only qualifying updates

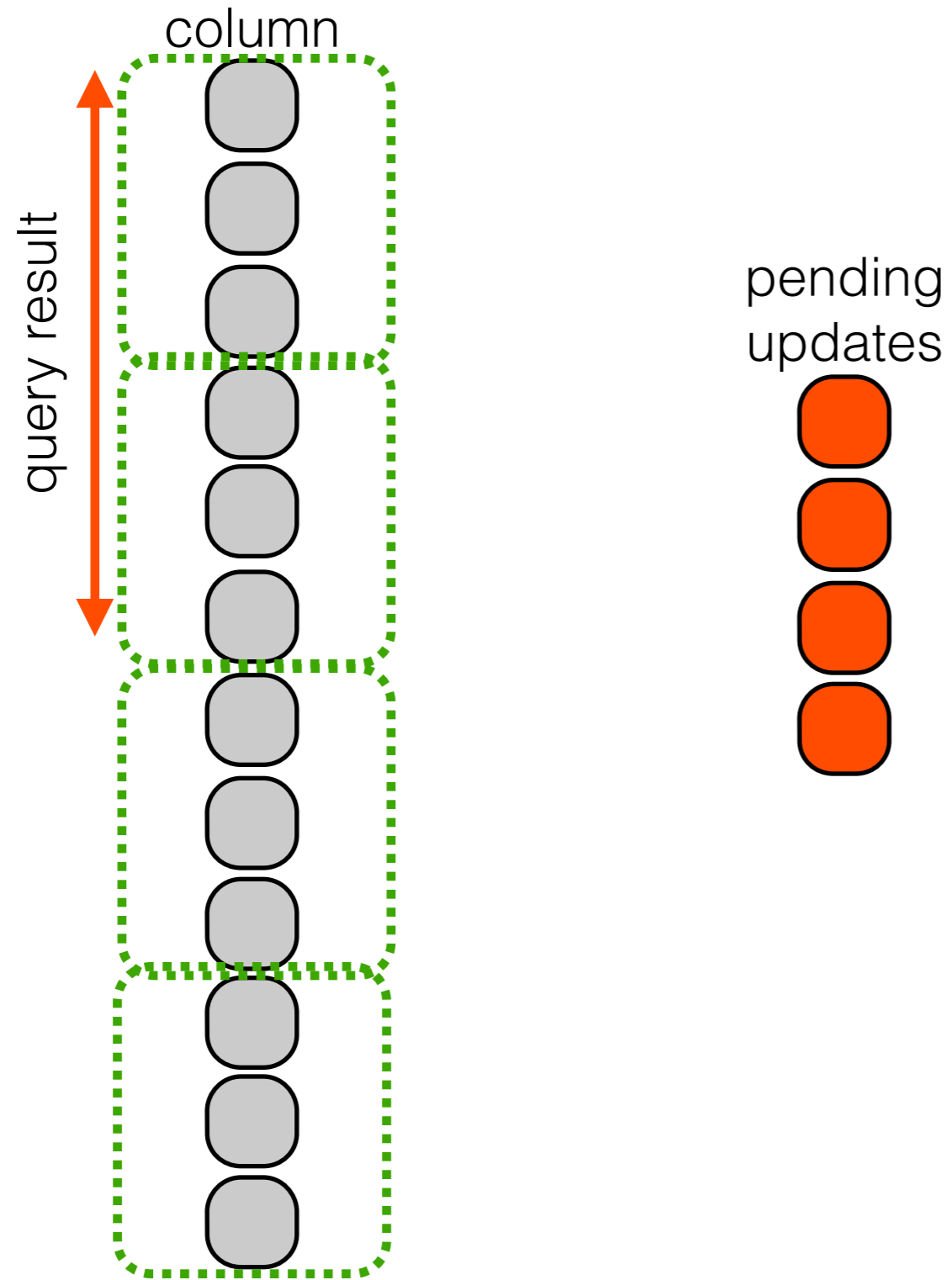


gradual merging avoids high picks

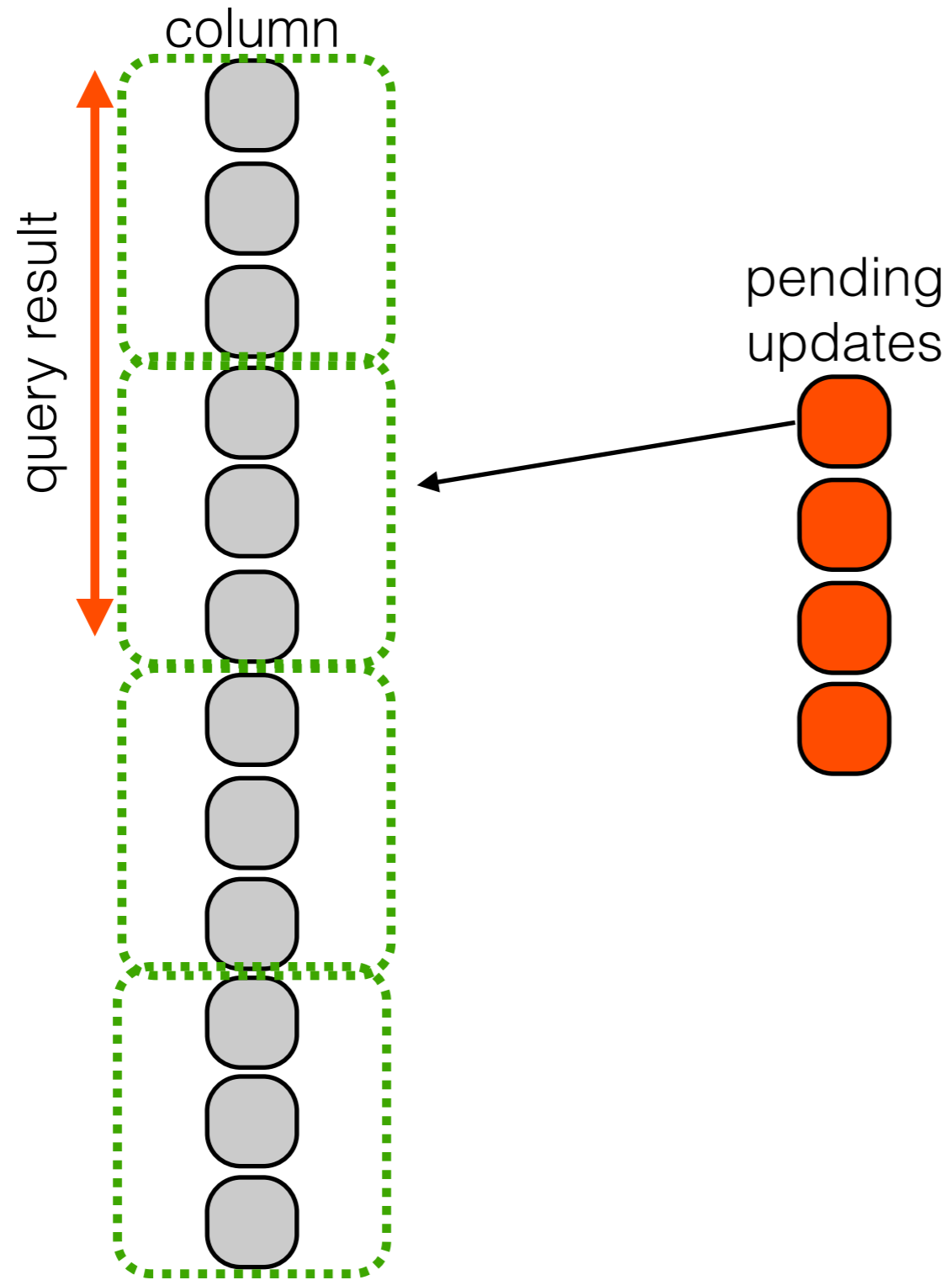
the ripple



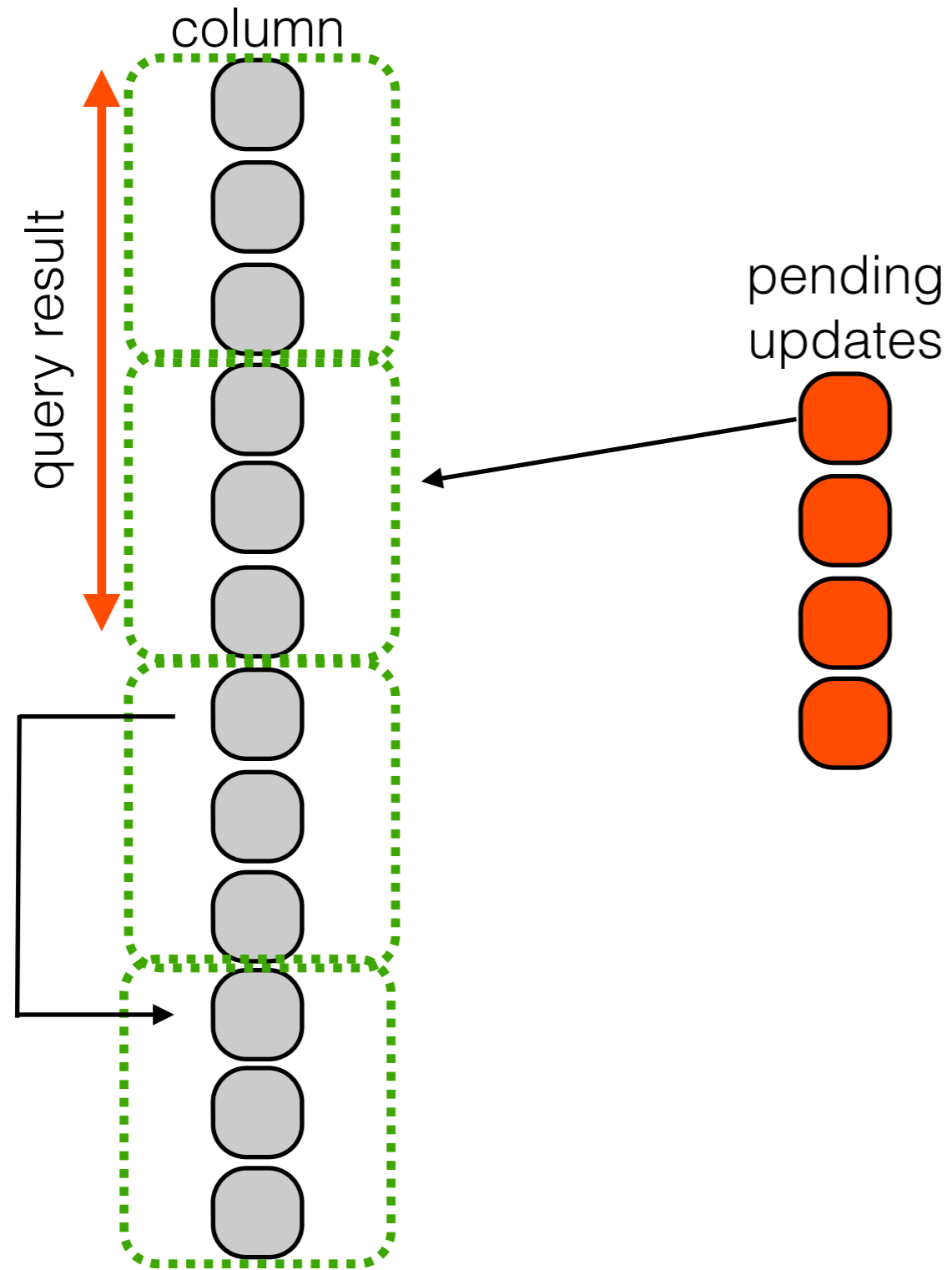
the ripple



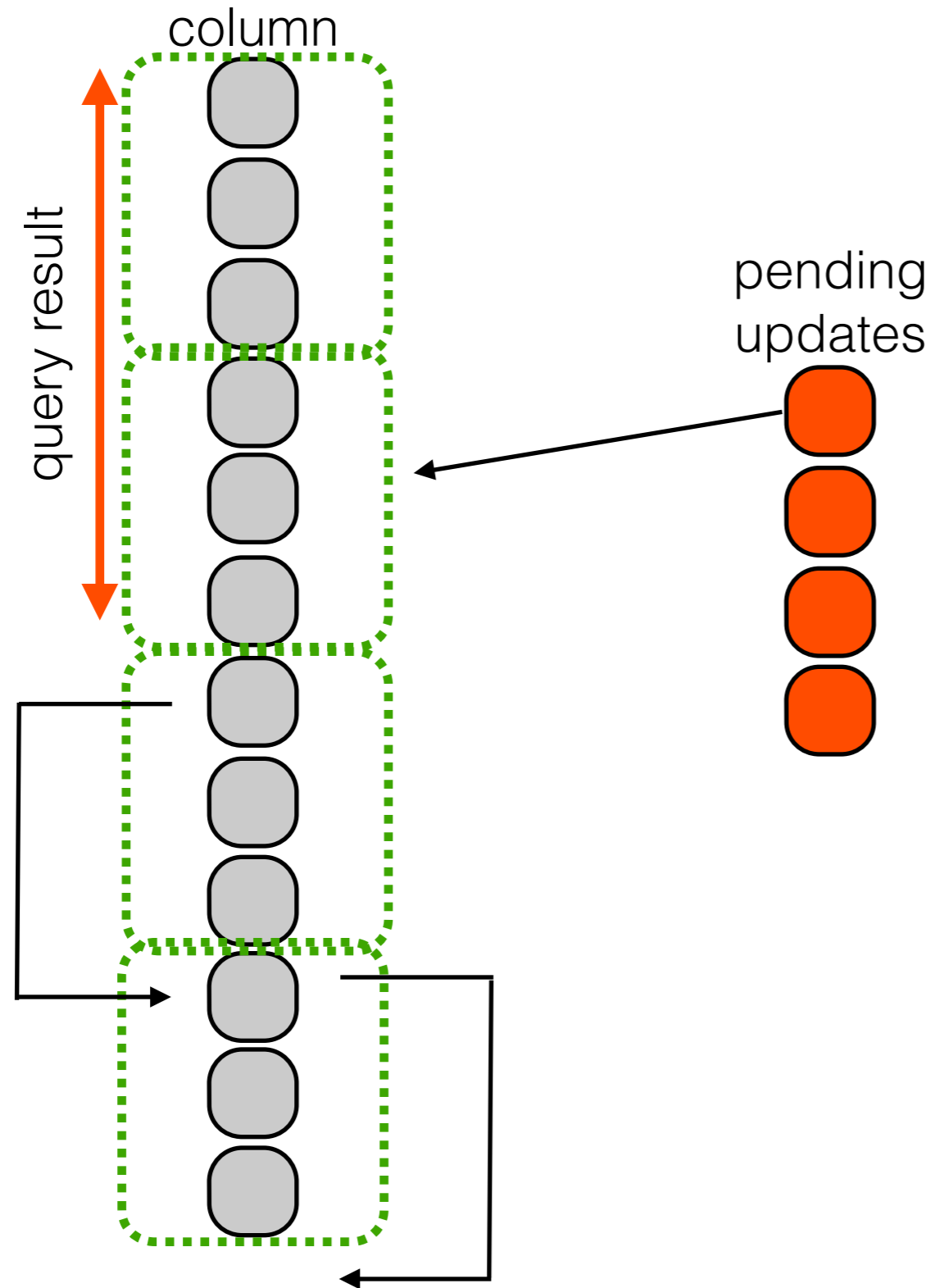
the ripple



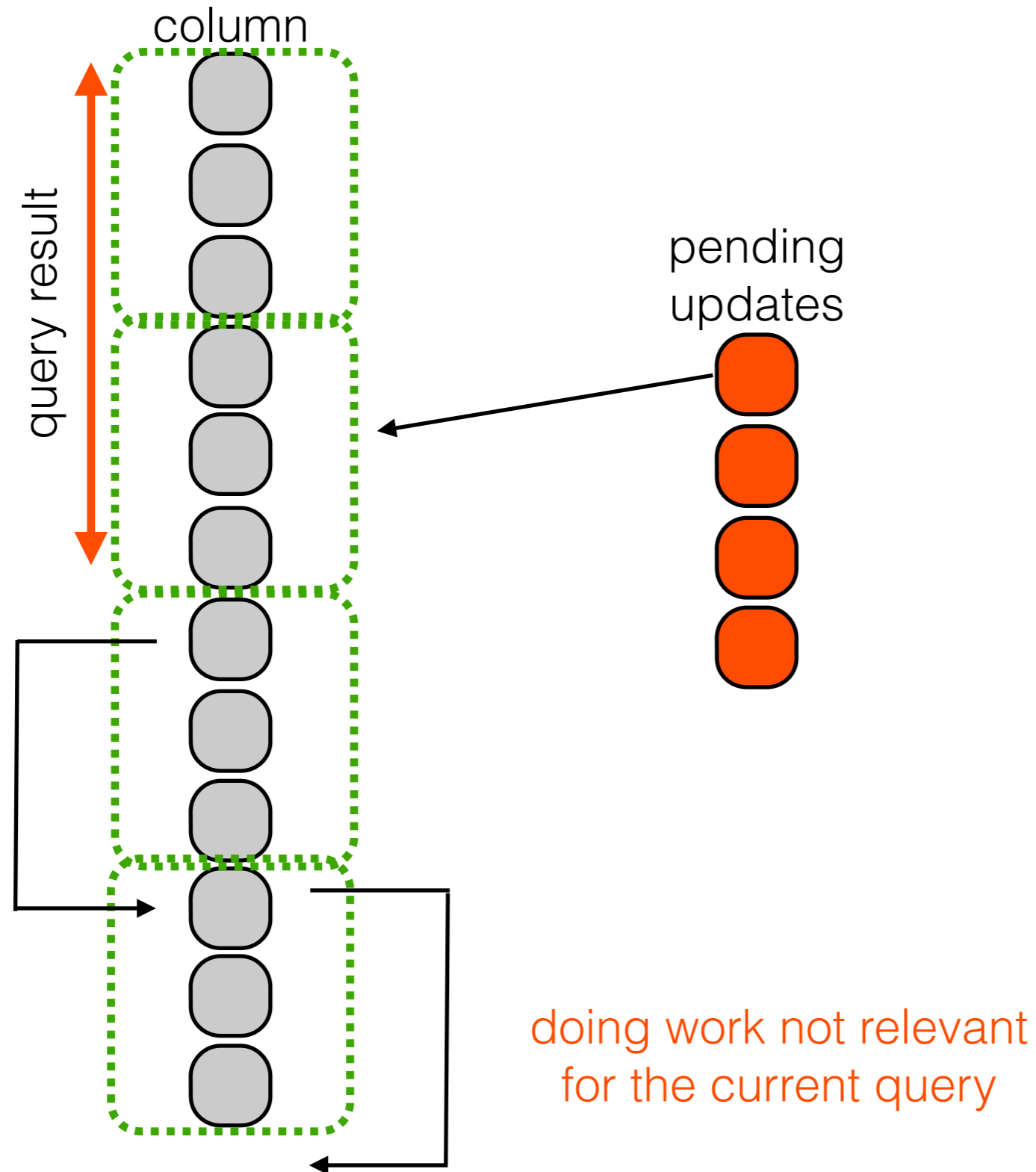
the ripple



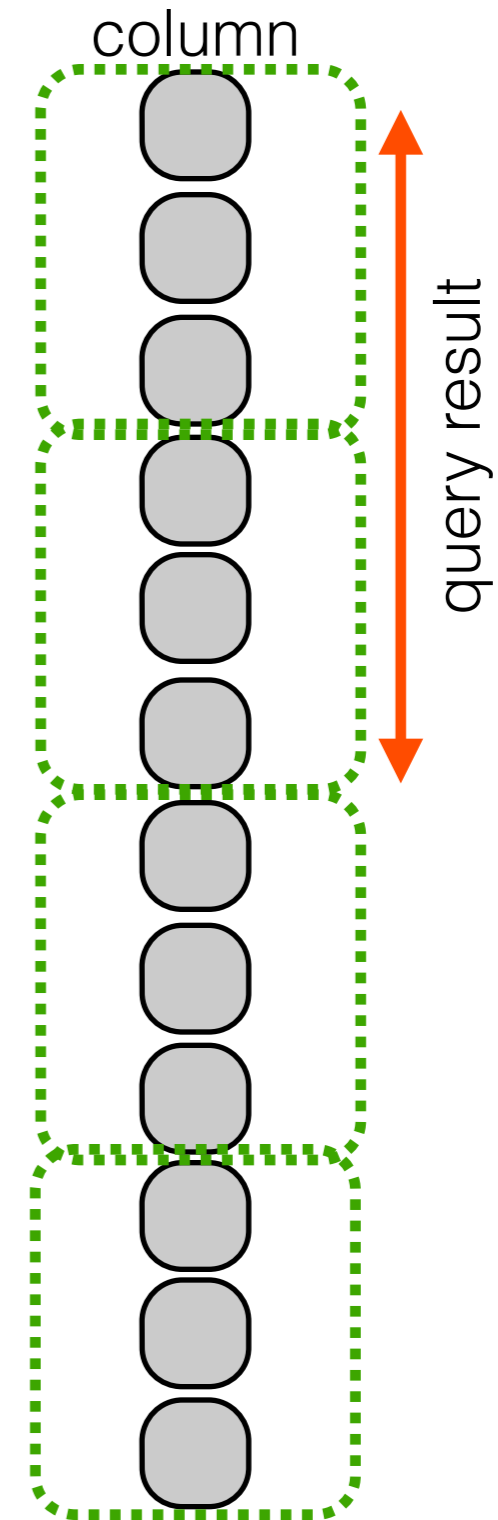
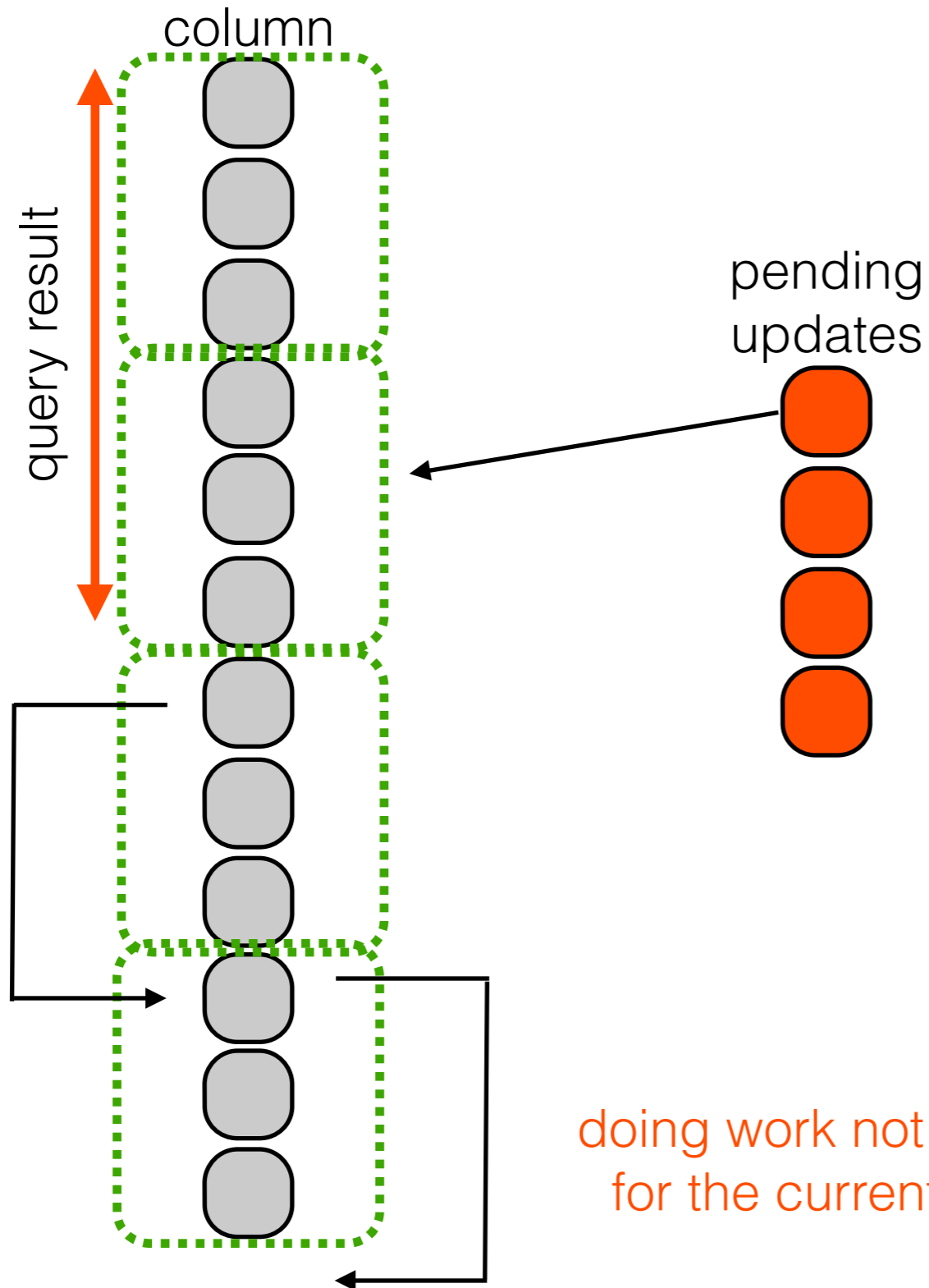
the ripple



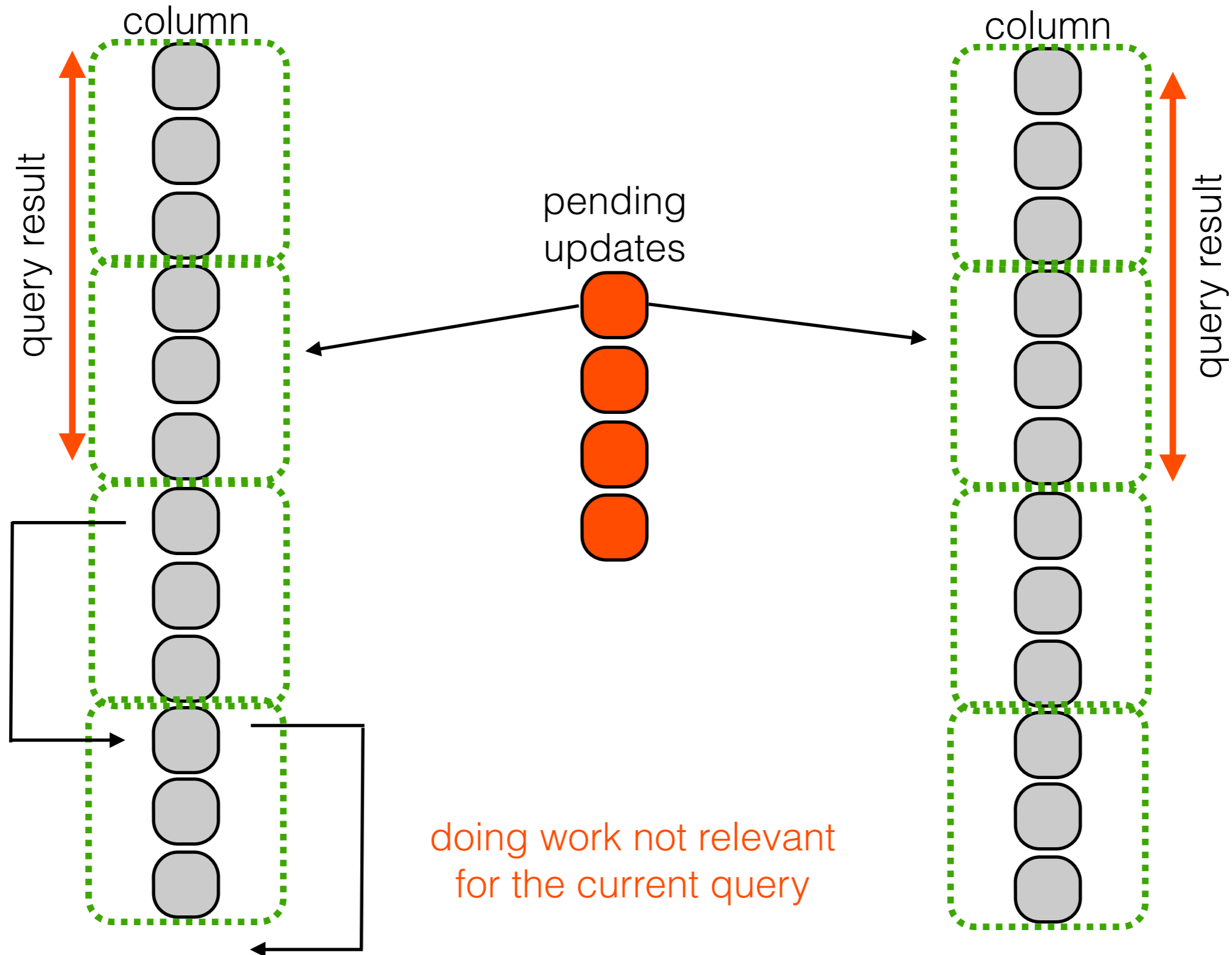
the ripple



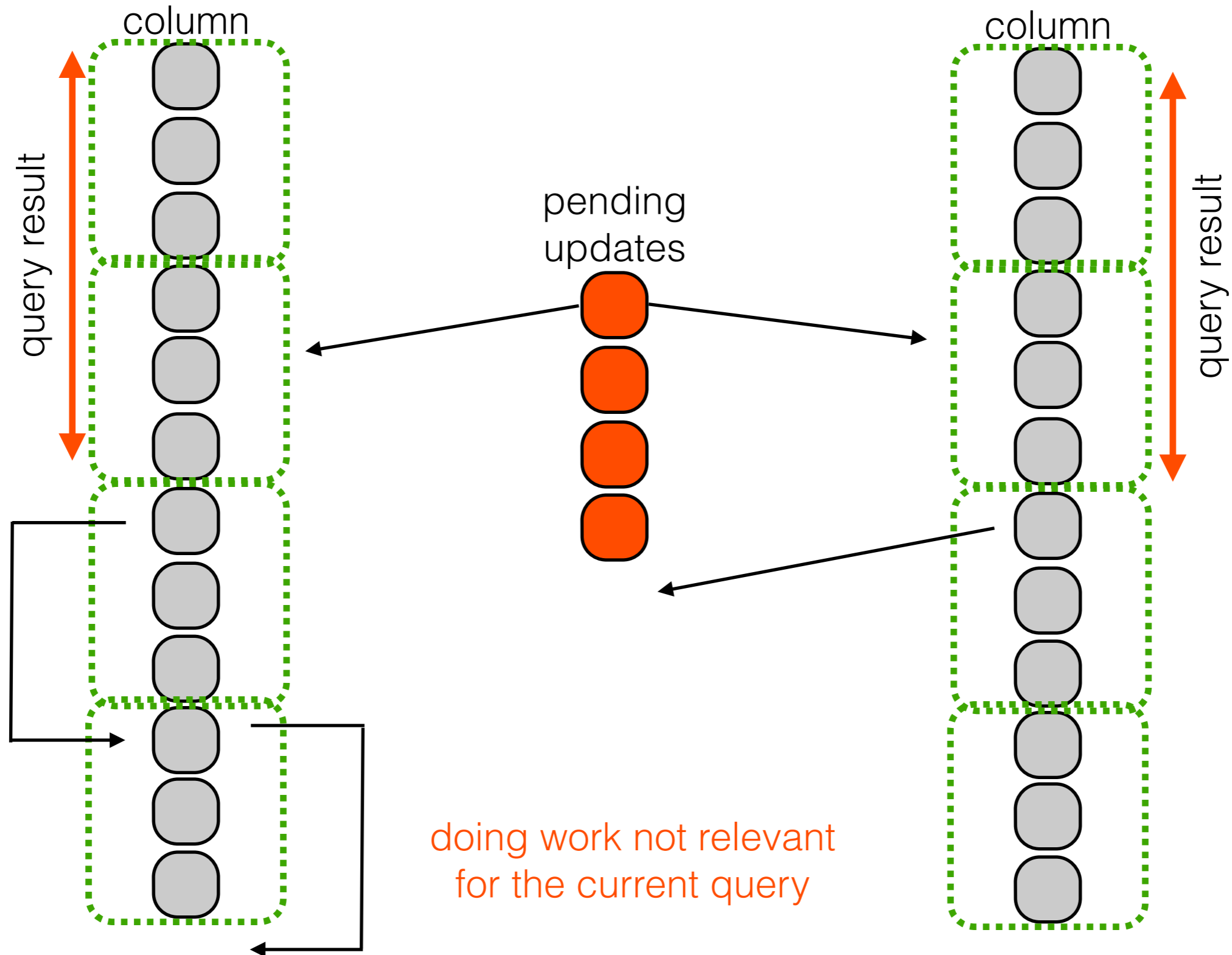
the ripple



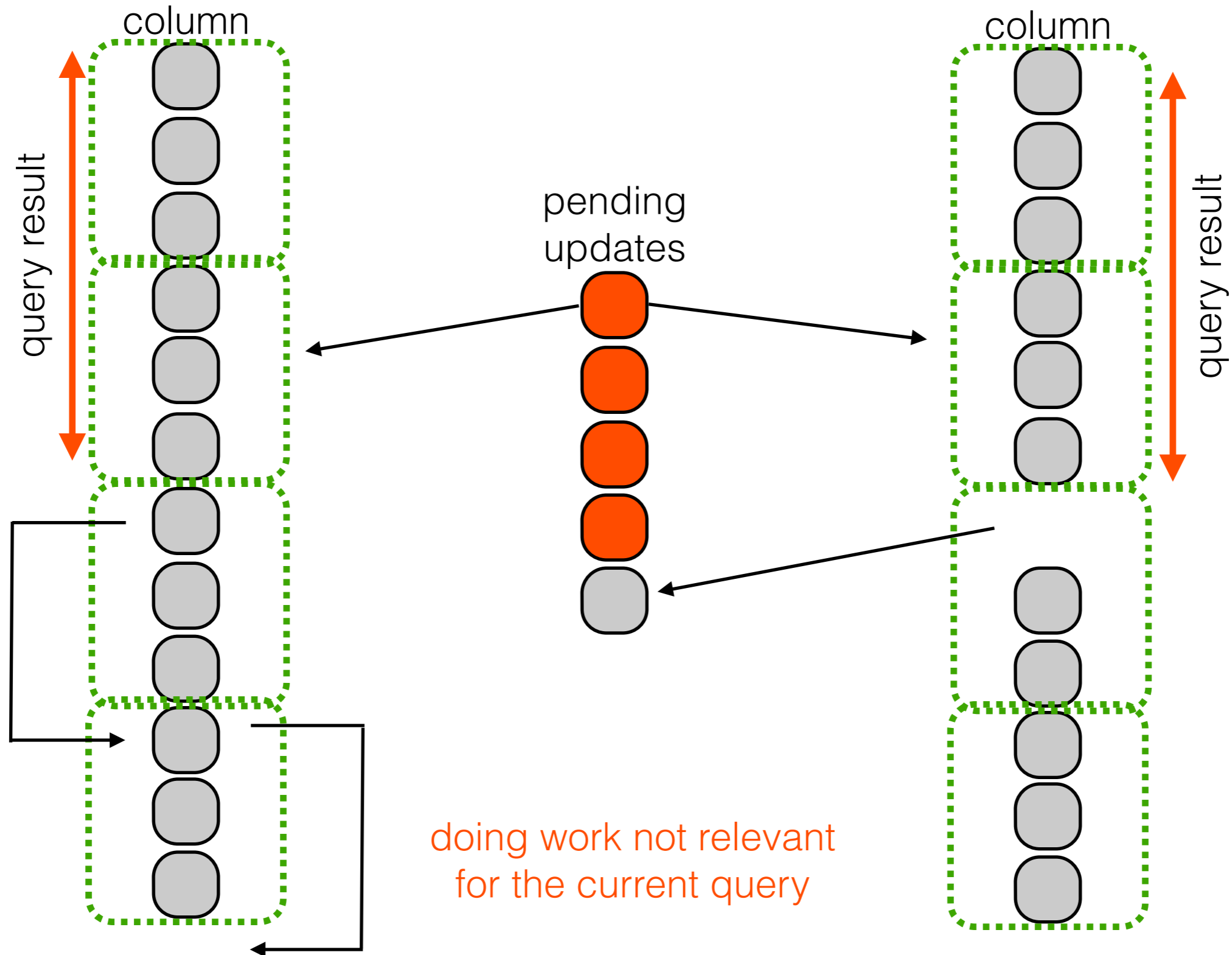
the ripple



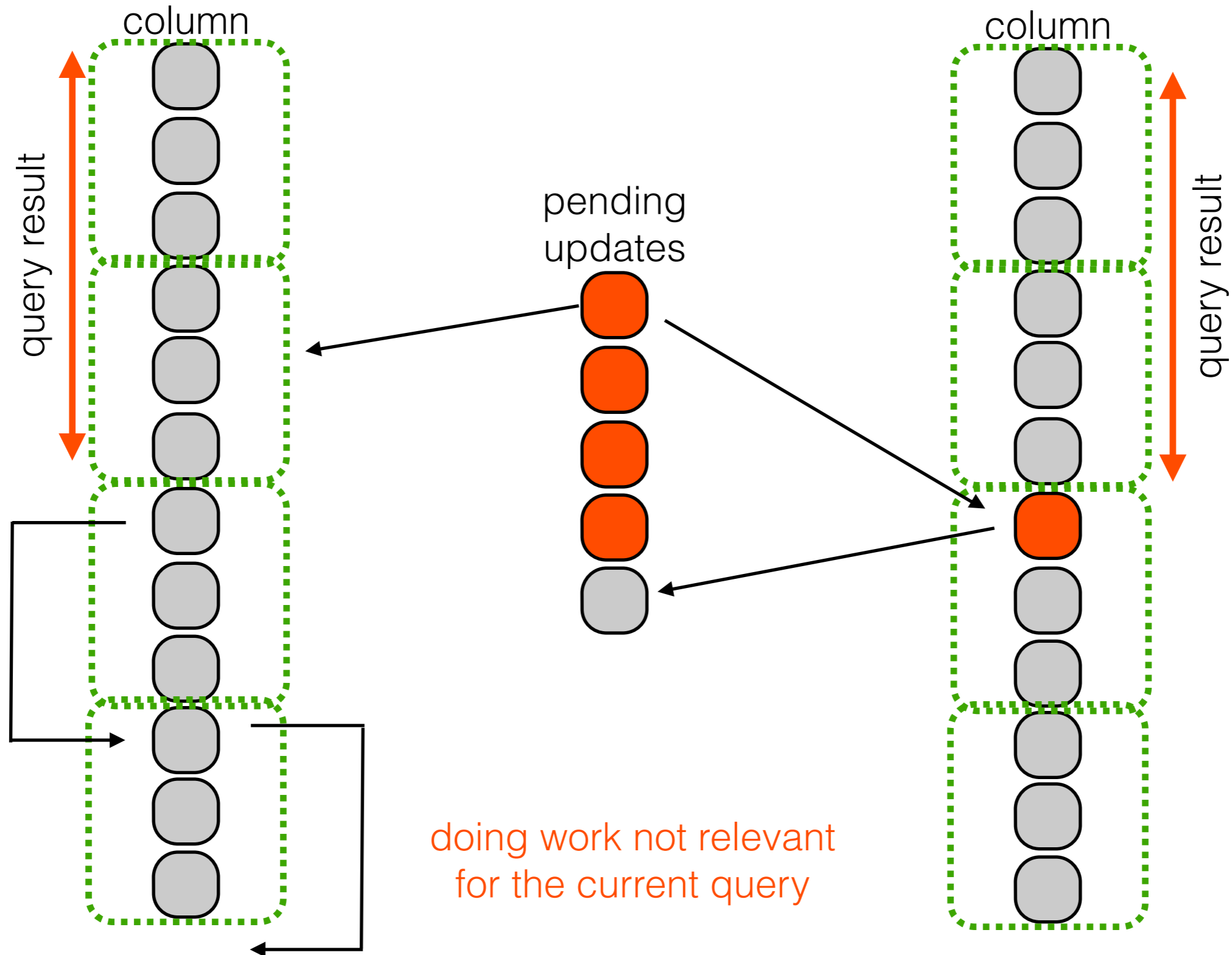
the ripple



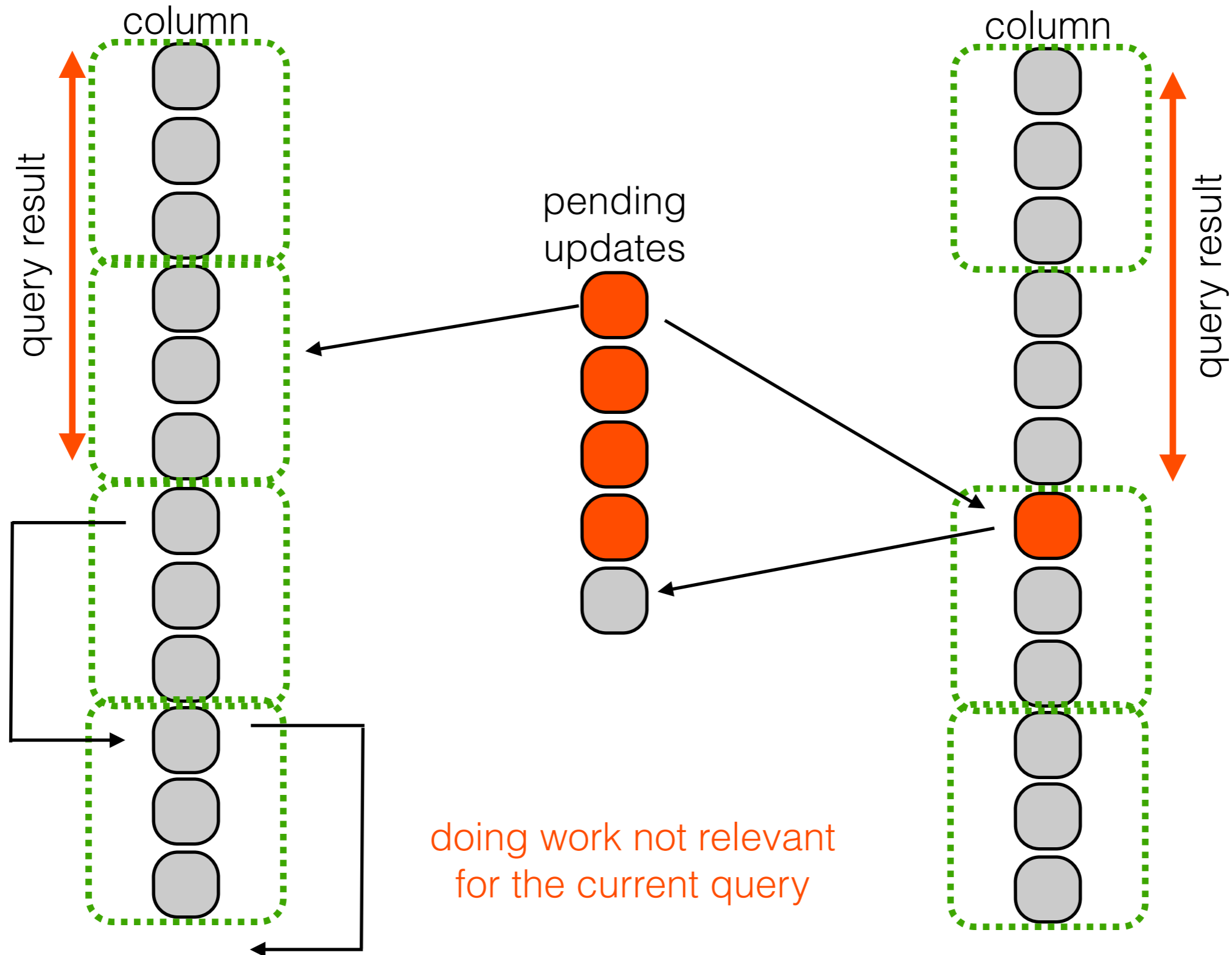
the ripple



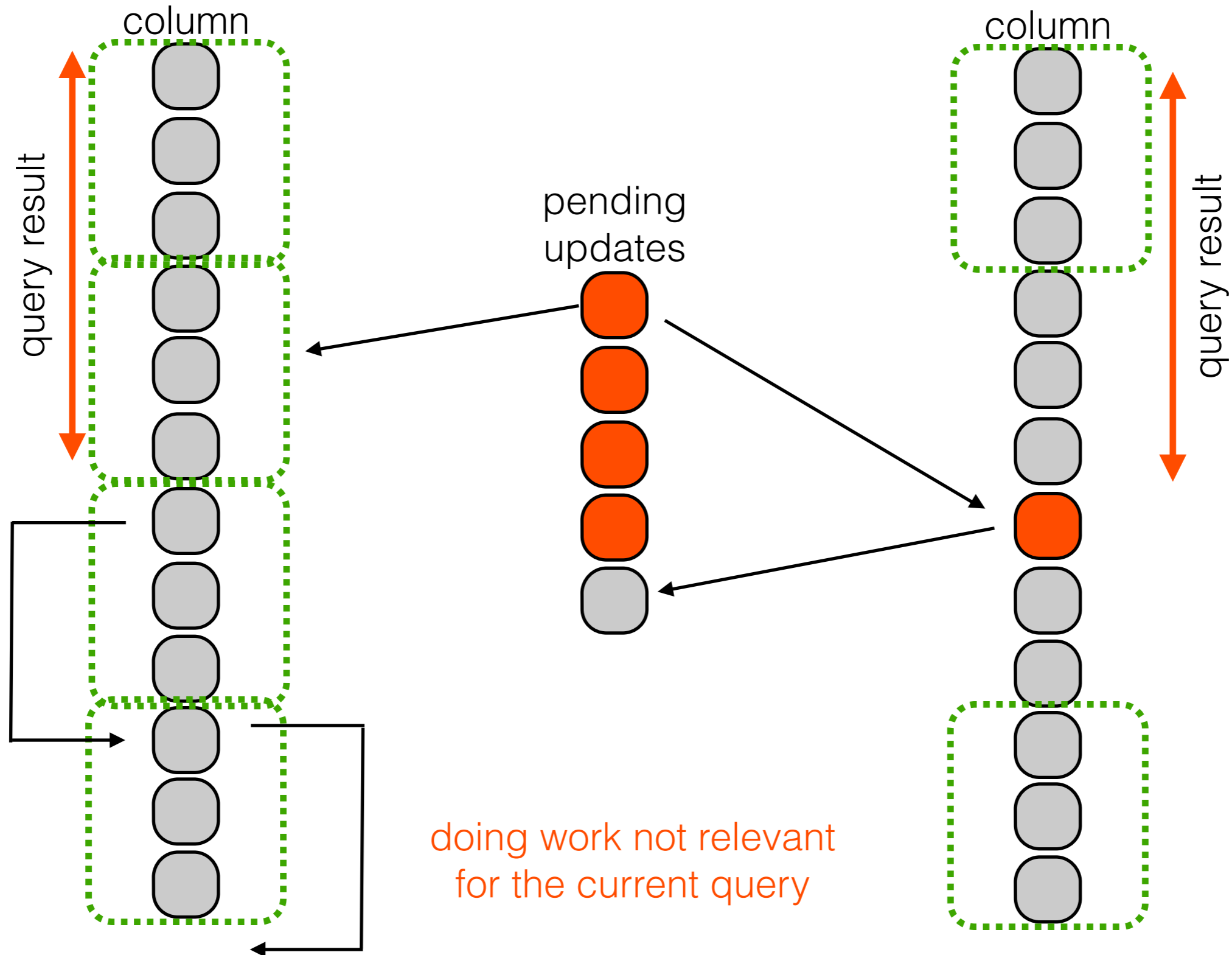
the ripple



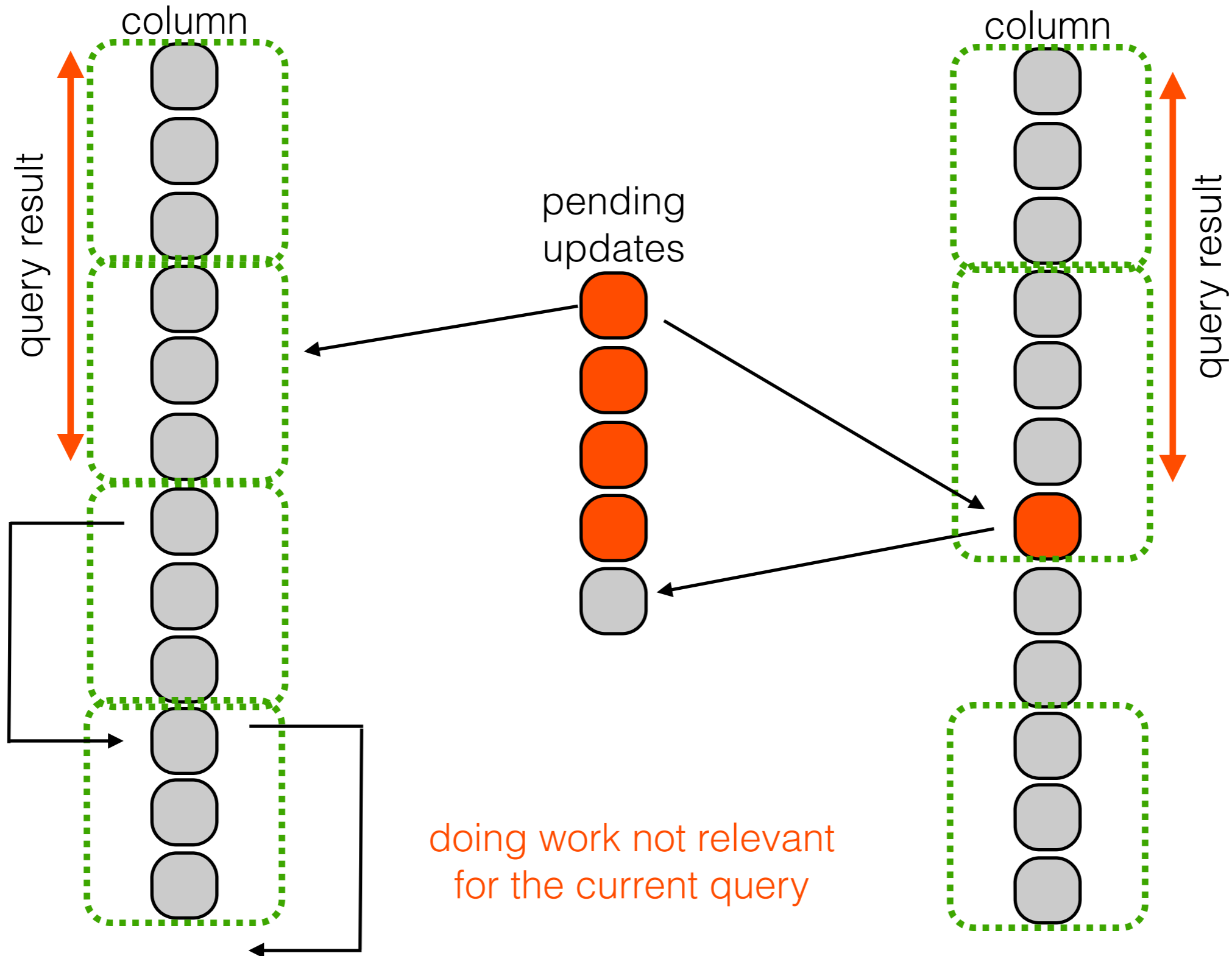
the ripple



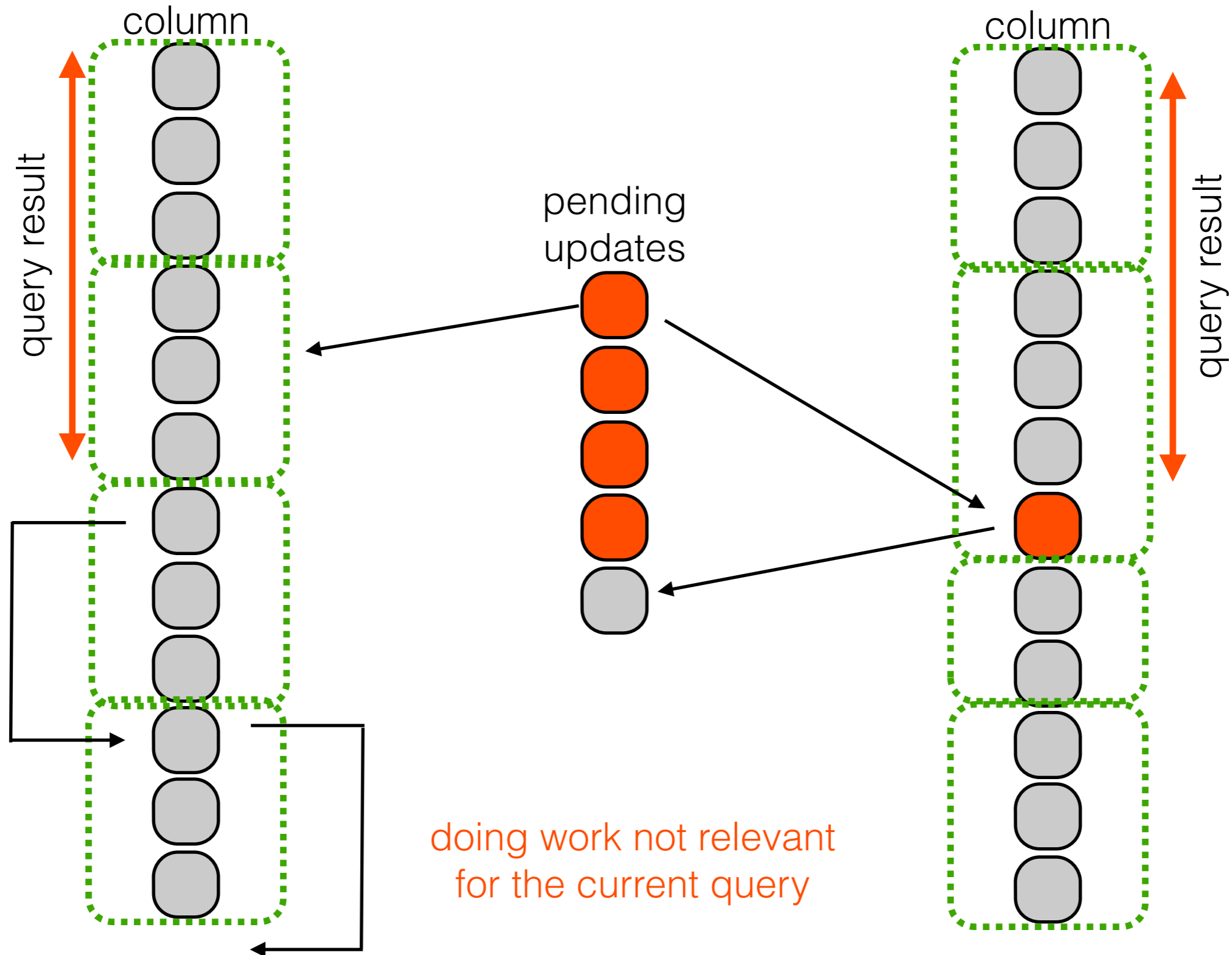
the ripple



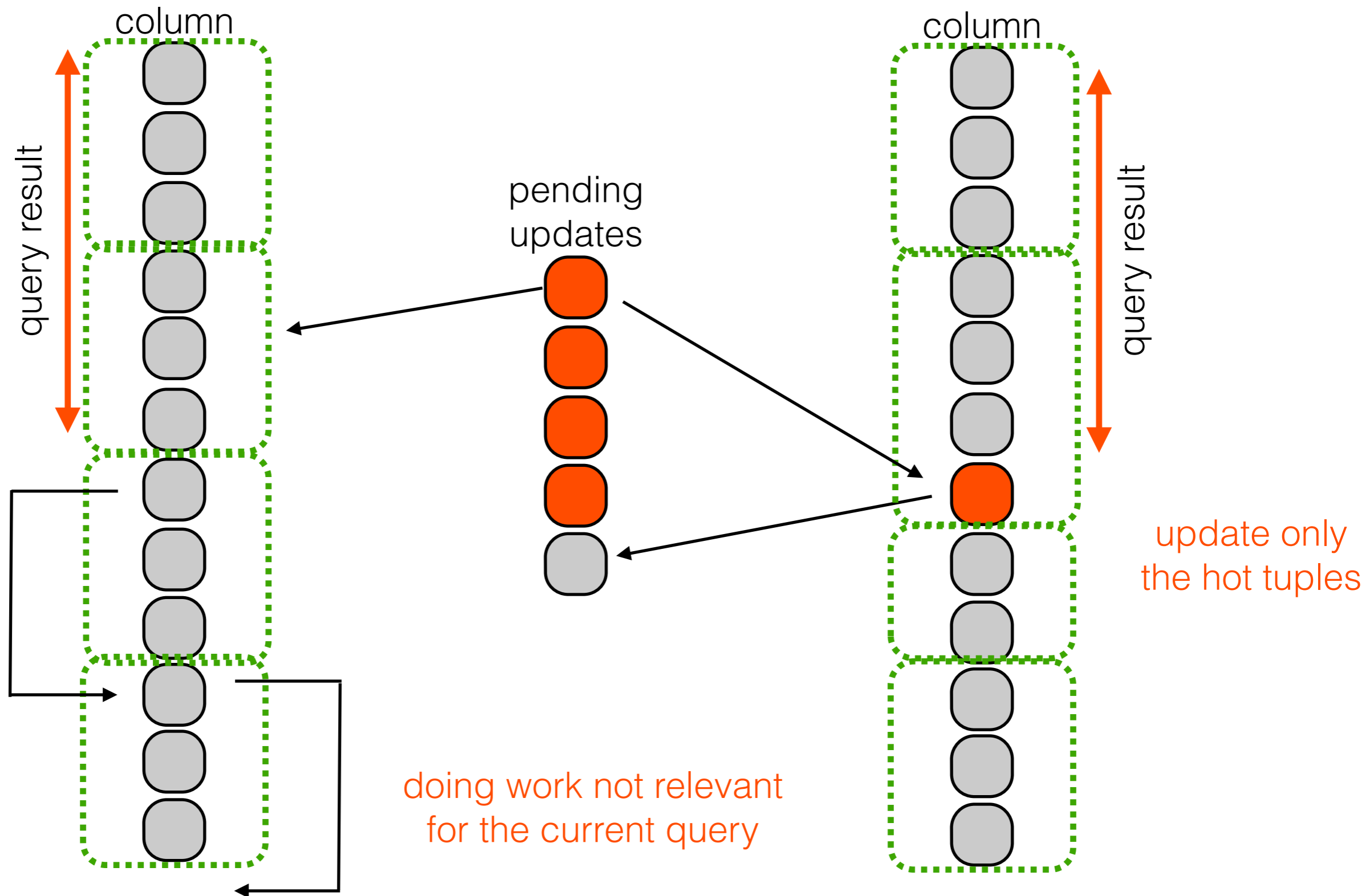
the ripple



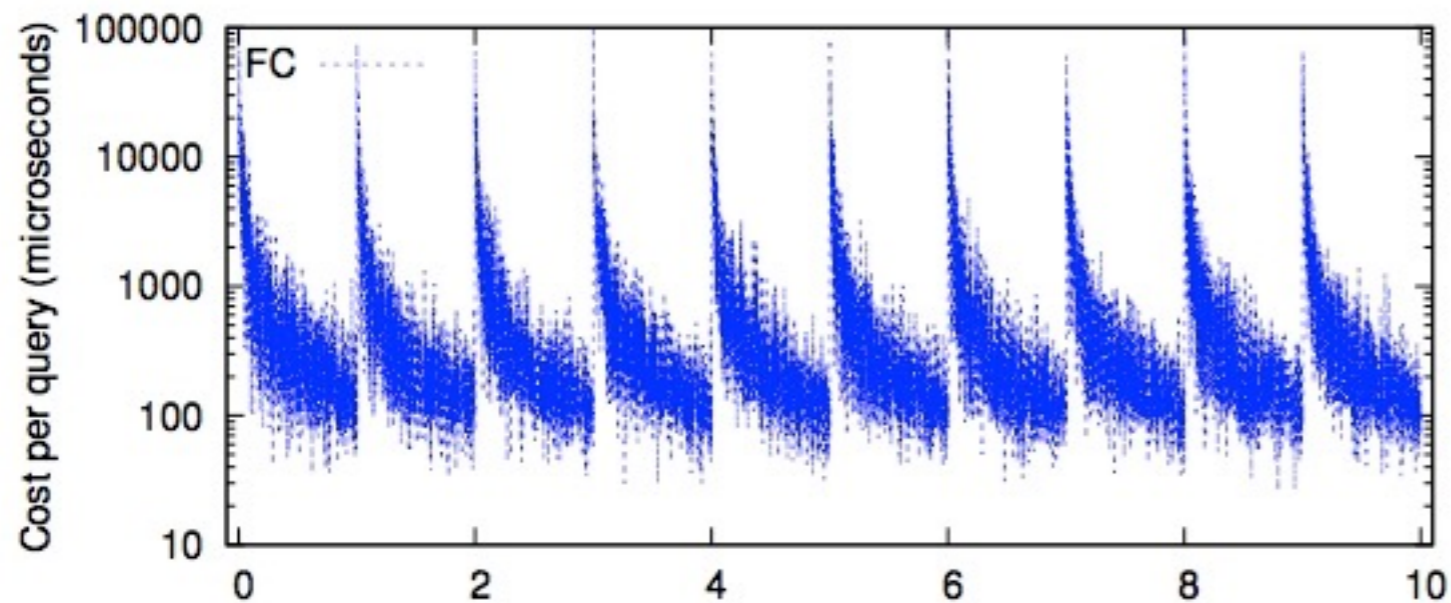
the ripple



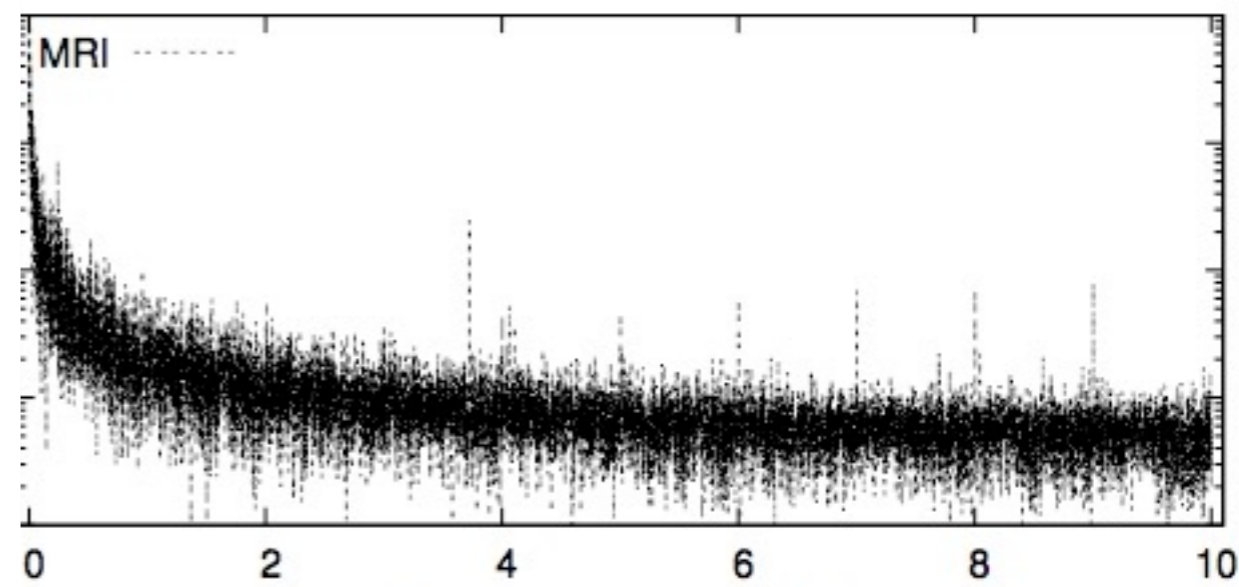
the ripple



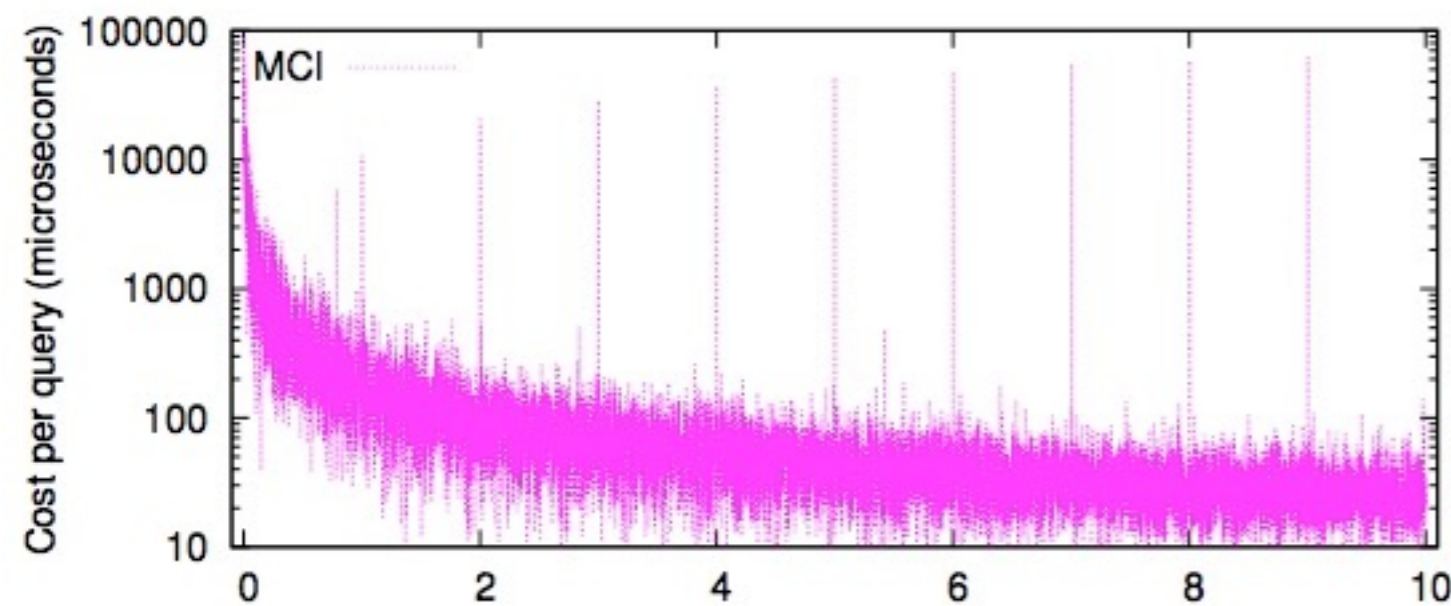
forget



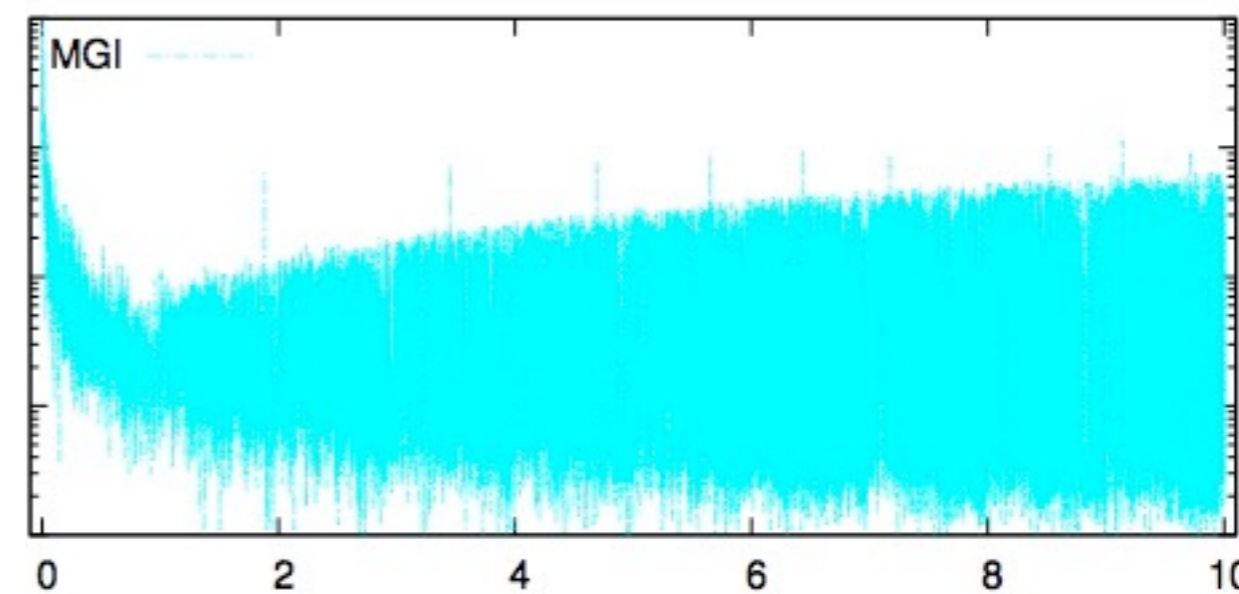
ripple



merge all updates



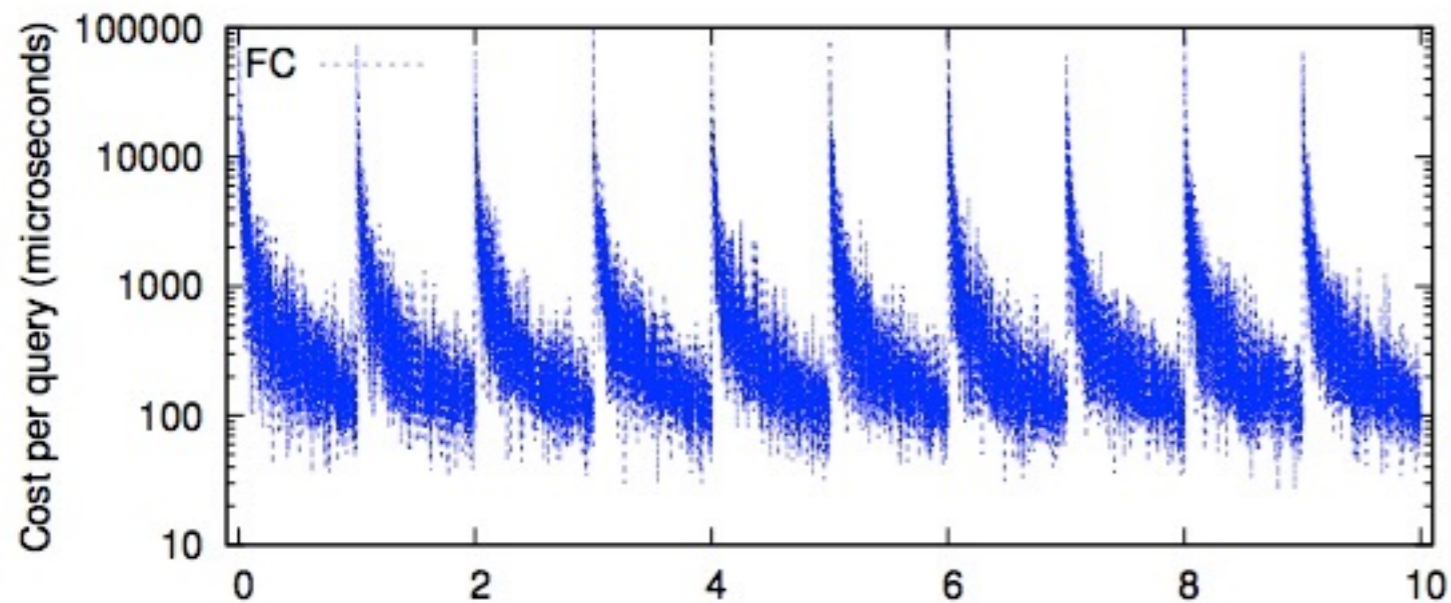
merge only qualifying updates



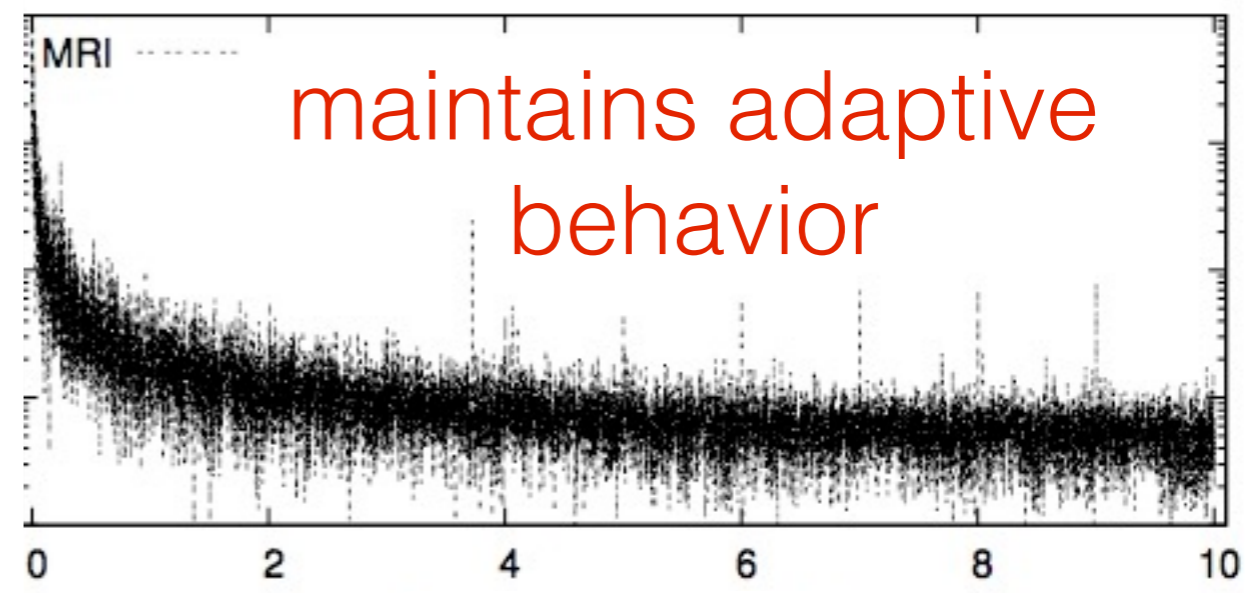
queries (x1000)

queries (x1000)

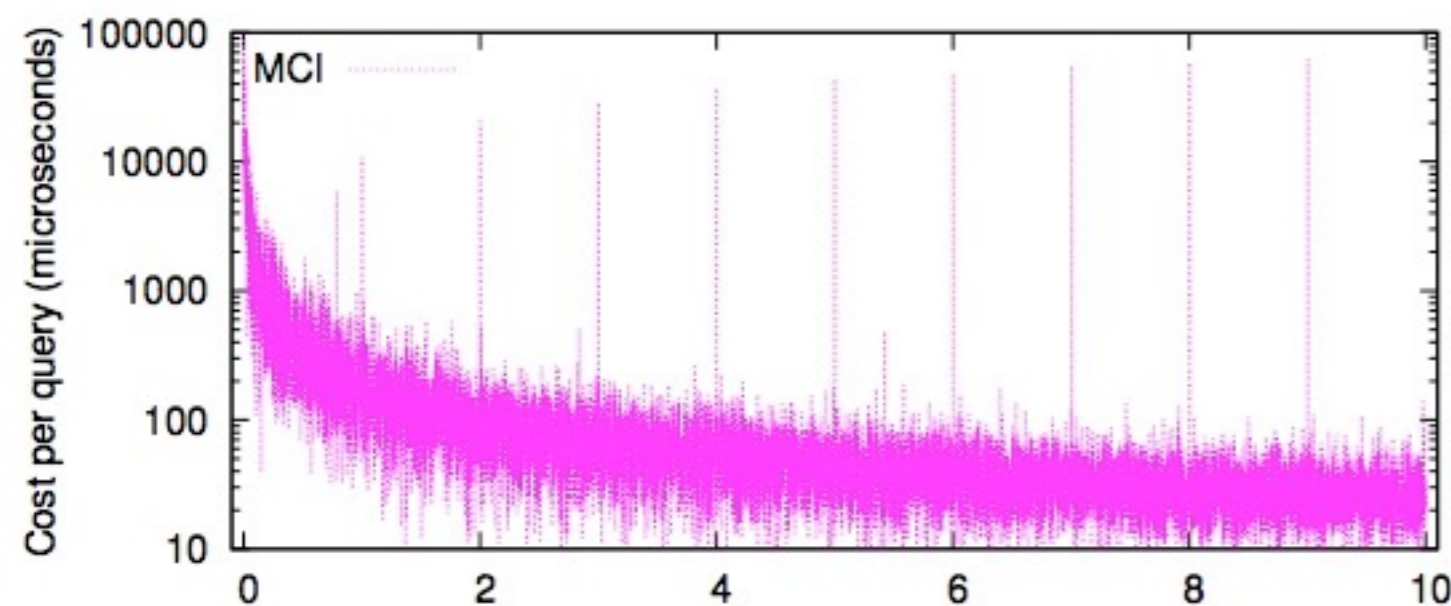
forget



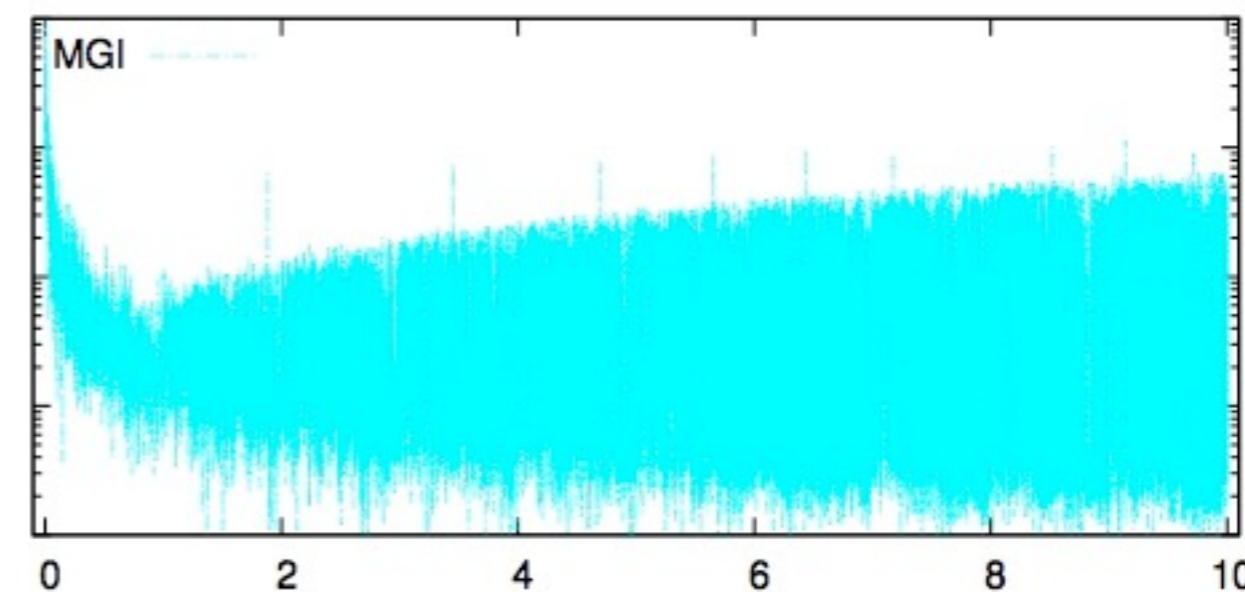
ripple



merge all updates



merge only qualifying updates



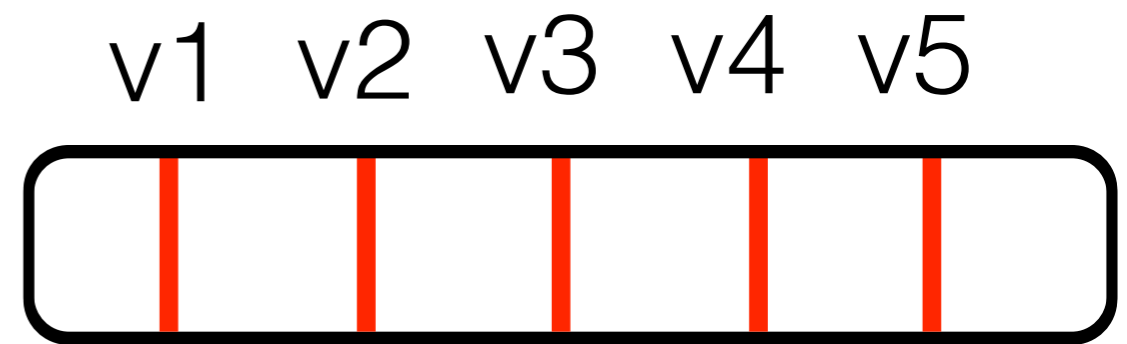
queries (x1000)

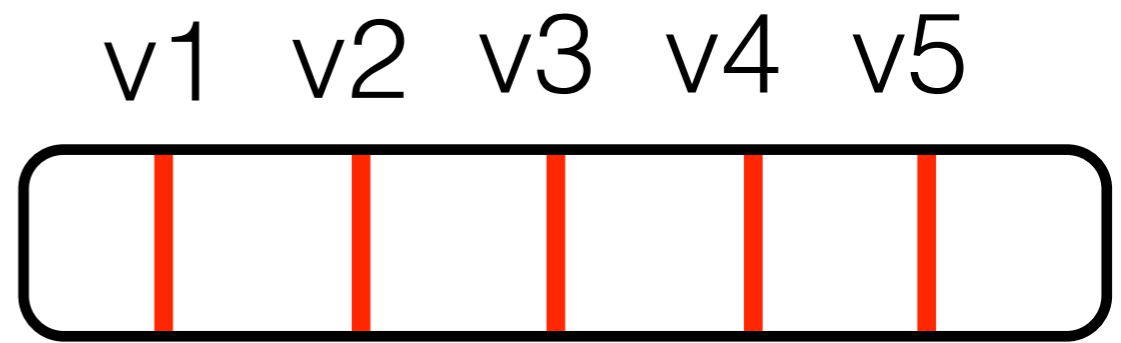
queries (x1000)

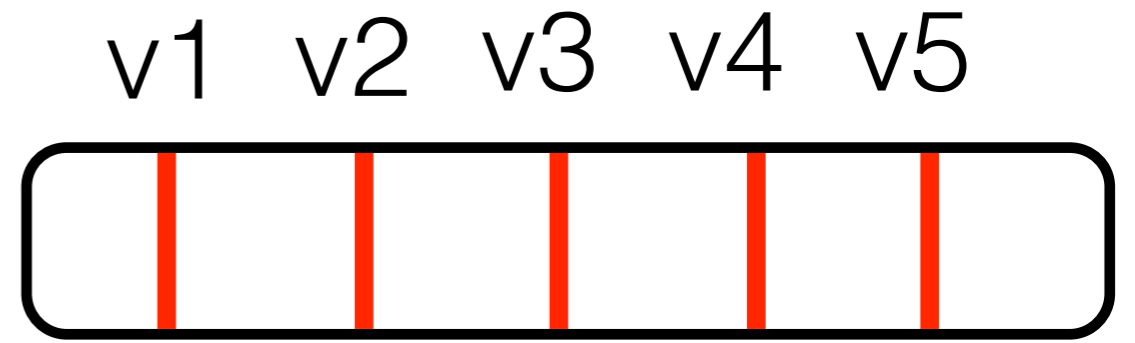
stochastic cracking

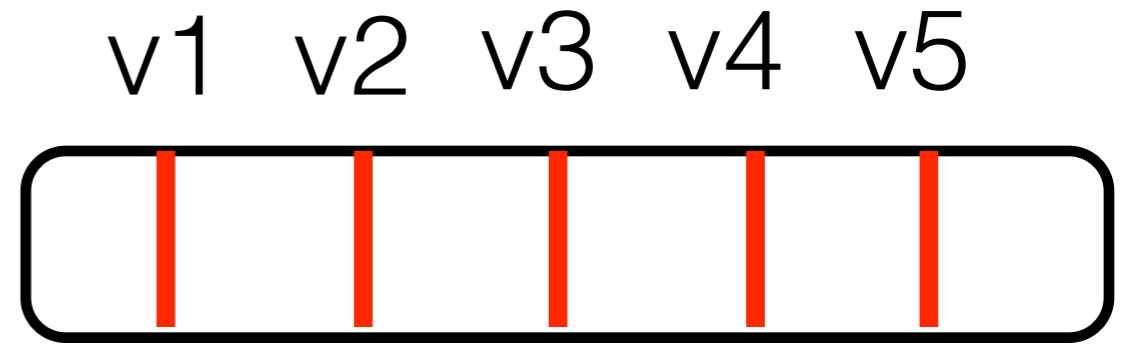
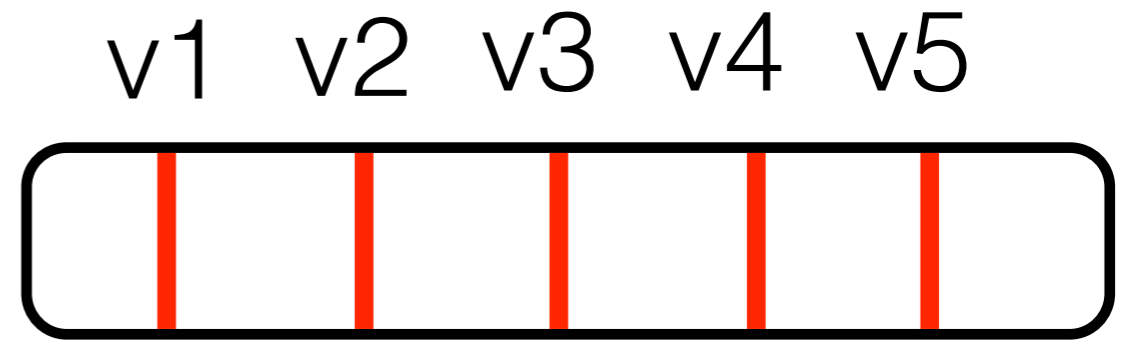
robustness











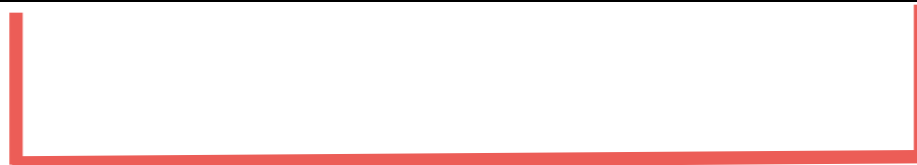




select [15,55]



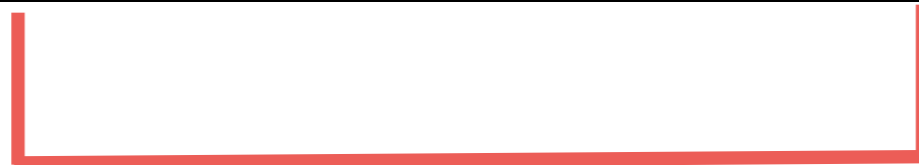
select [15,55]



select [15,55]

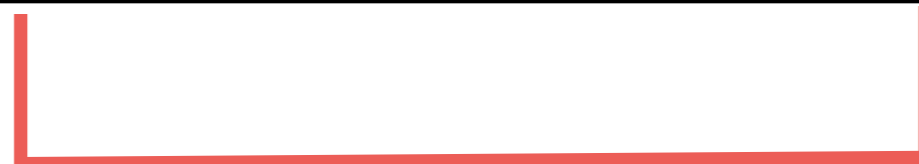
10 20 30 40 50 60



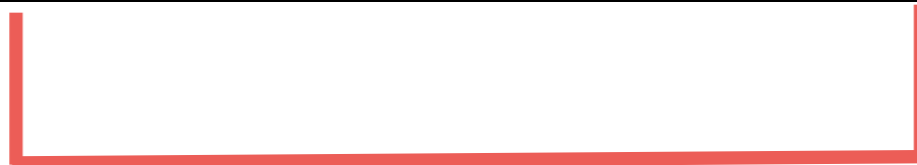


select [15,55]

10 20 30 40 50 60

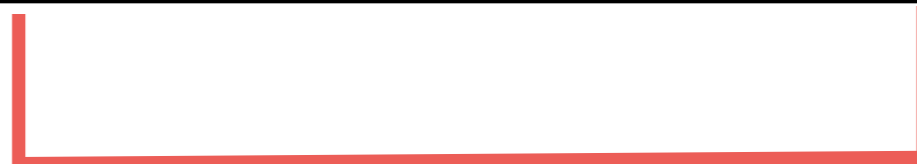


select [15,55]



select [15,55]

10 20 30 40 50 60



select [15,55]



select [15,55]

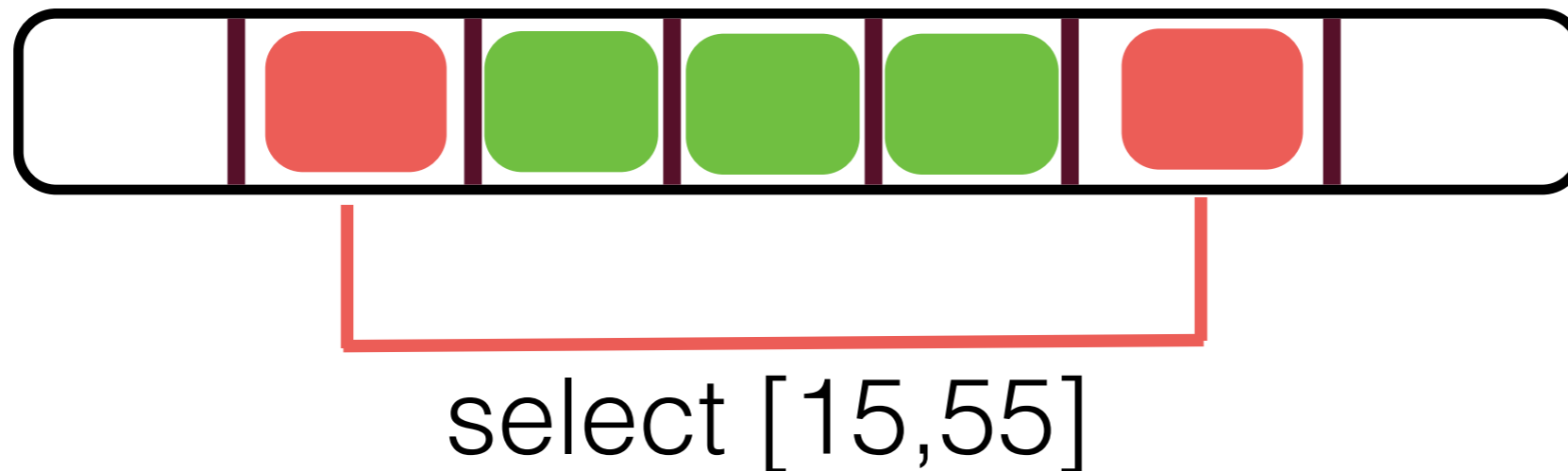
10 20 30 40 50 60



select [15,55]

the amount of work for each query depends on the index state

the state of the index depends on past queries patterns



column with 100 unique integers [1,100]



good sequence

column with 100 unique integers [1,100]

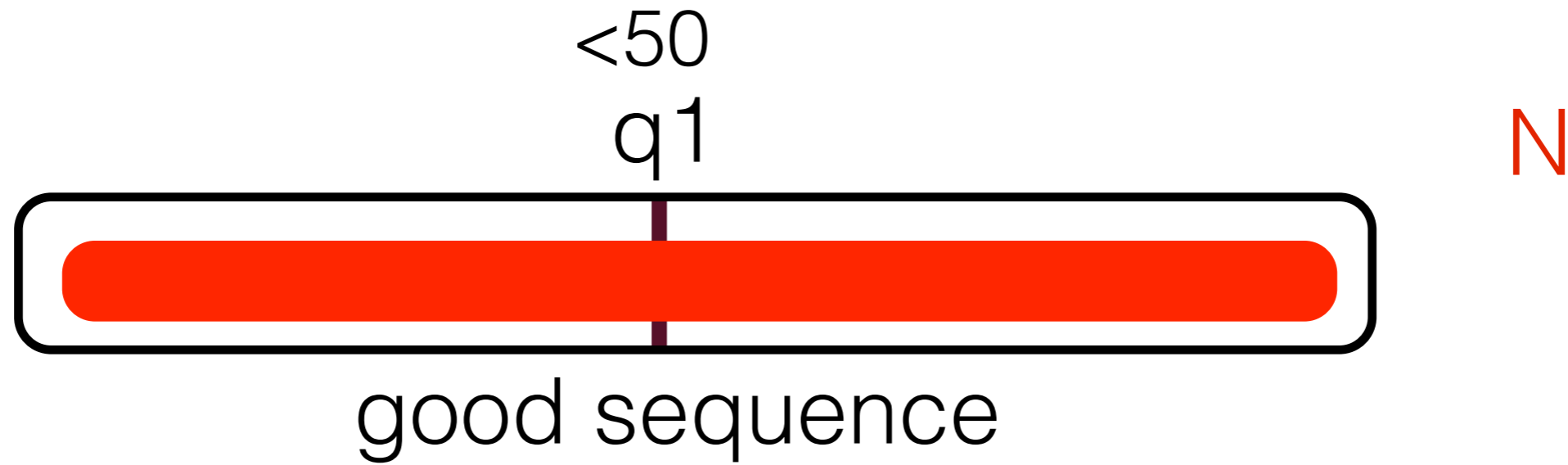
<50

q1



good sequence

column with 100 unique integers [1,100]



column with 100 unique integers [1,100]

<50

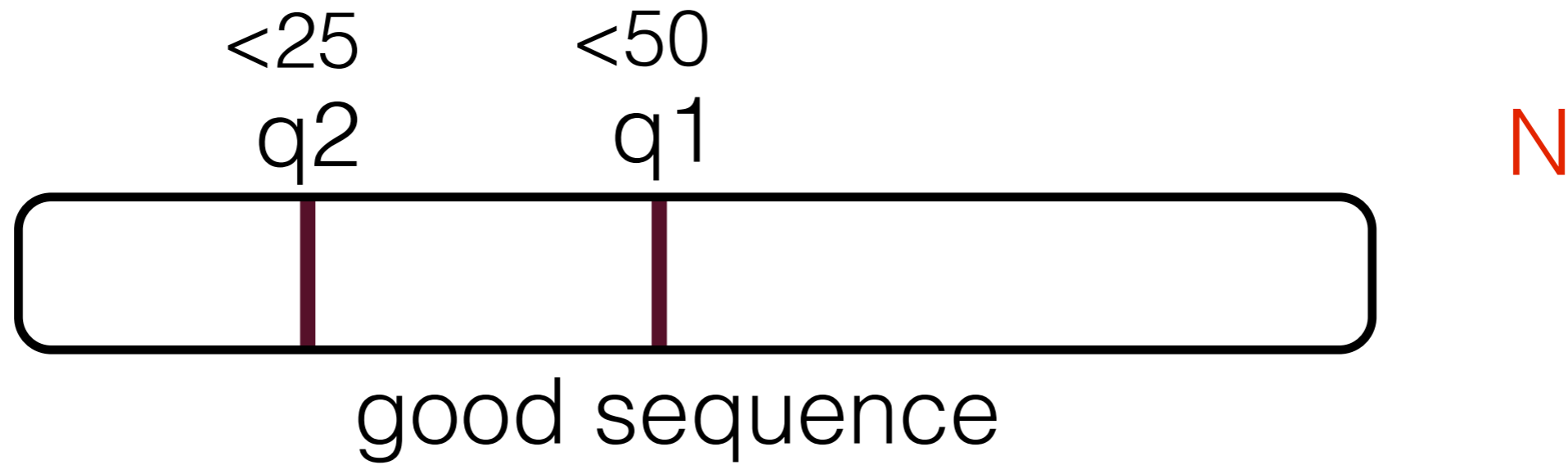
q1

N

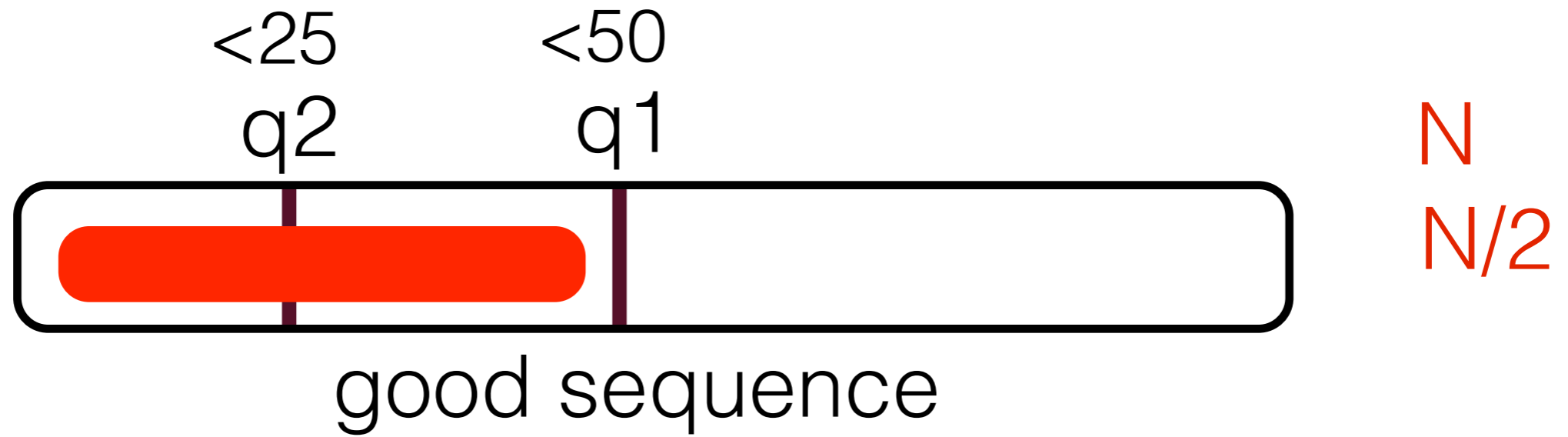


good sequence

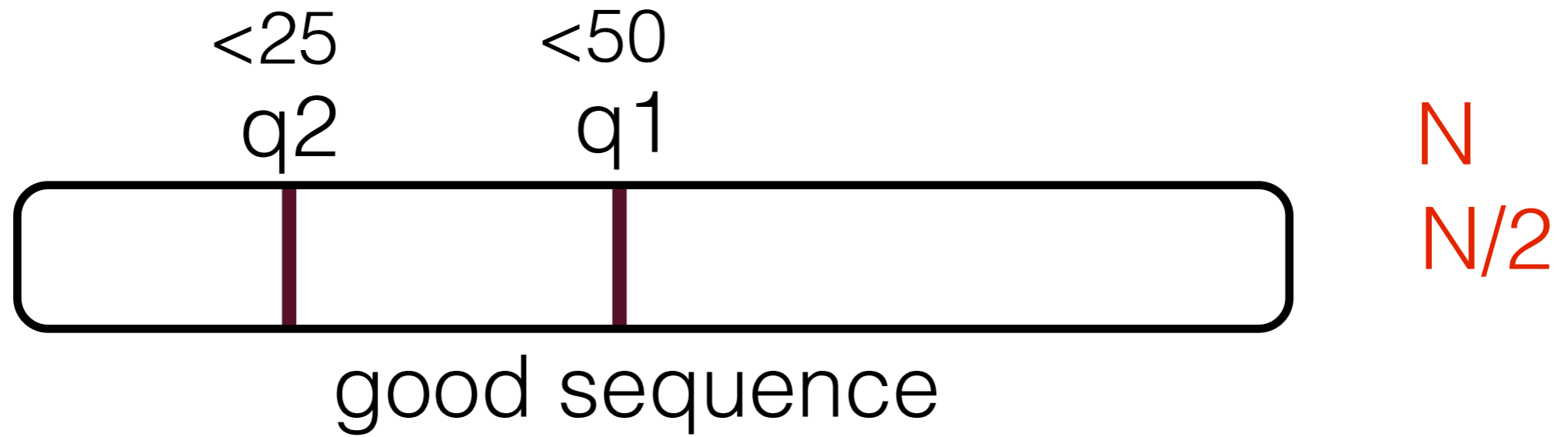
column with 100 unique integers [1,100]



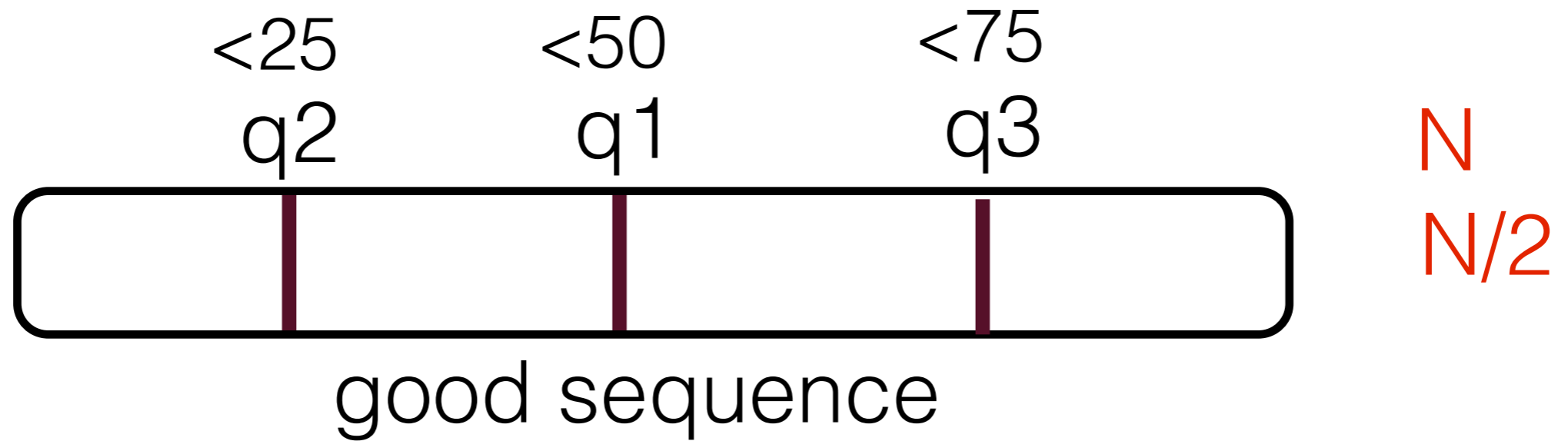
column with 100 unique integers [1,100]



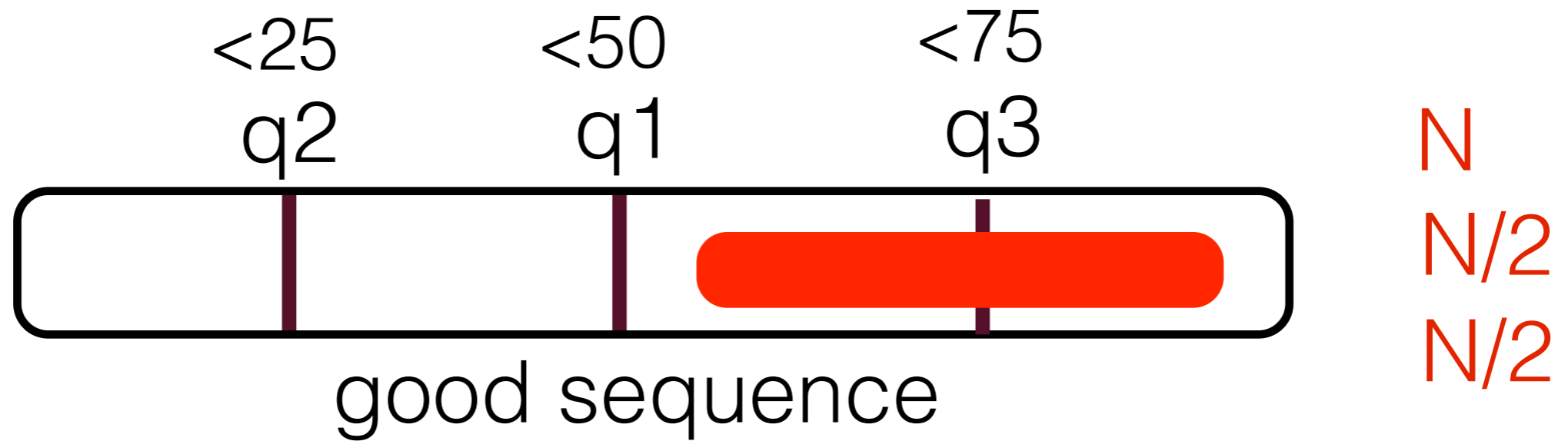
column with 100 unique integers [1,100]



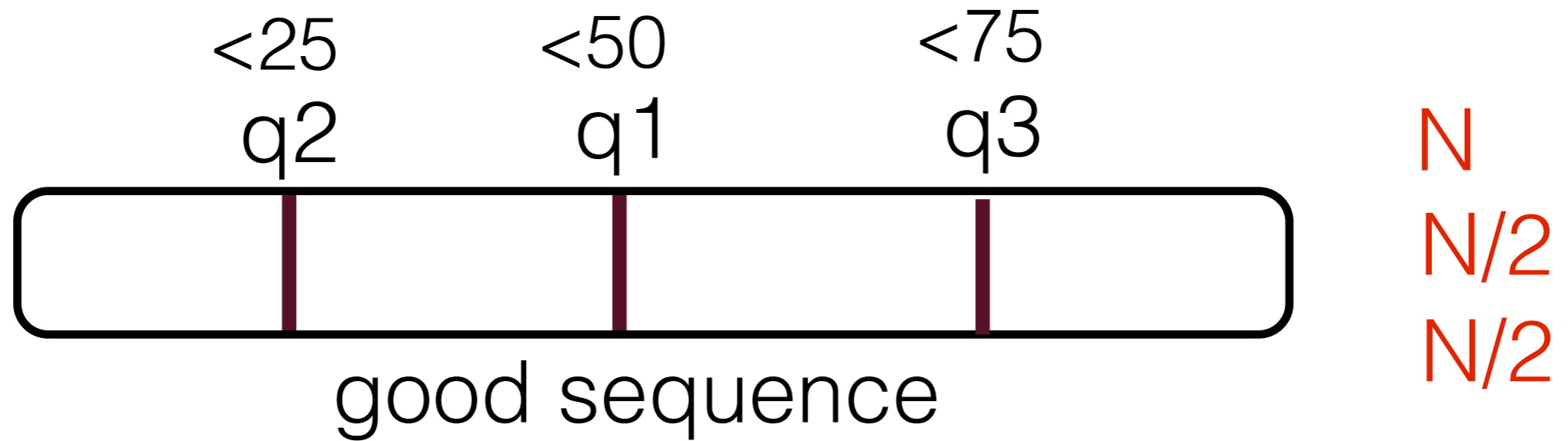
column with 100 unique integers [1,100]



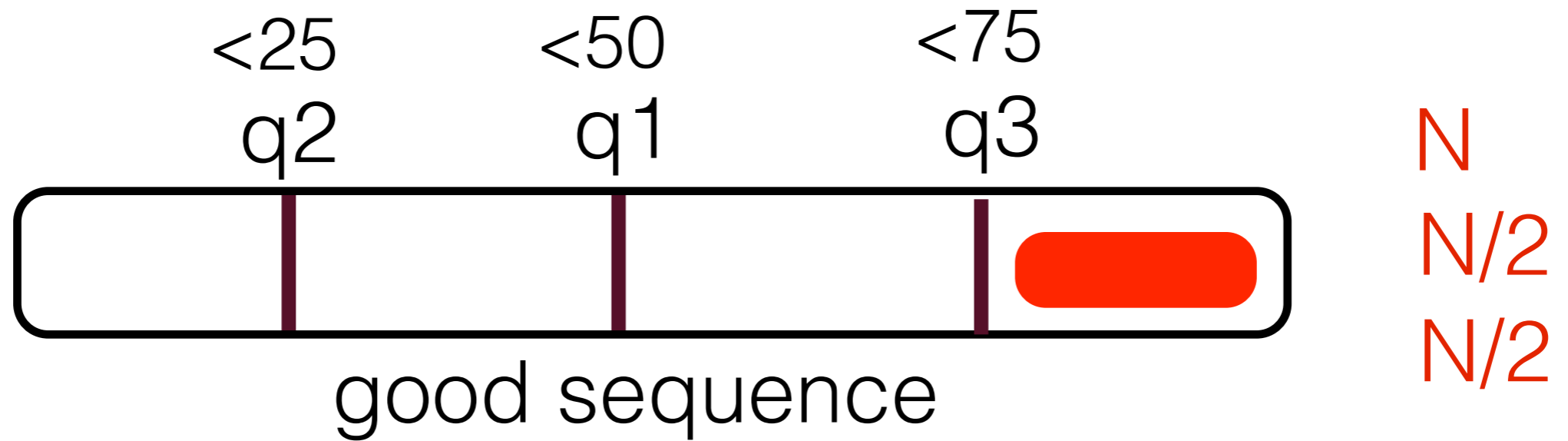
column with 100 unique integers [1,100]



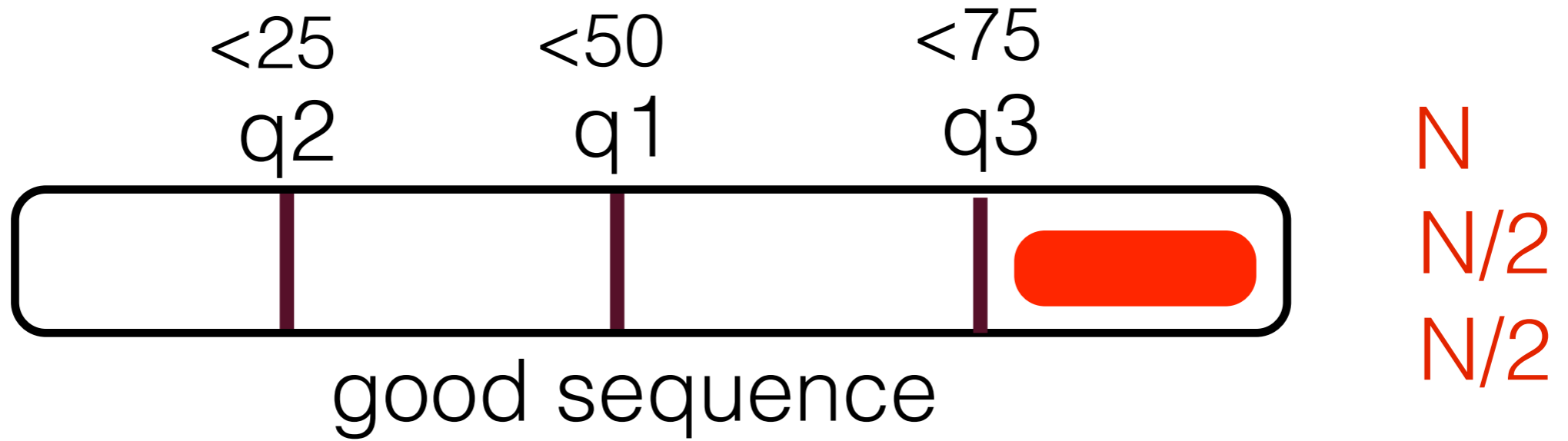
column with 100 unique integers [1,100]



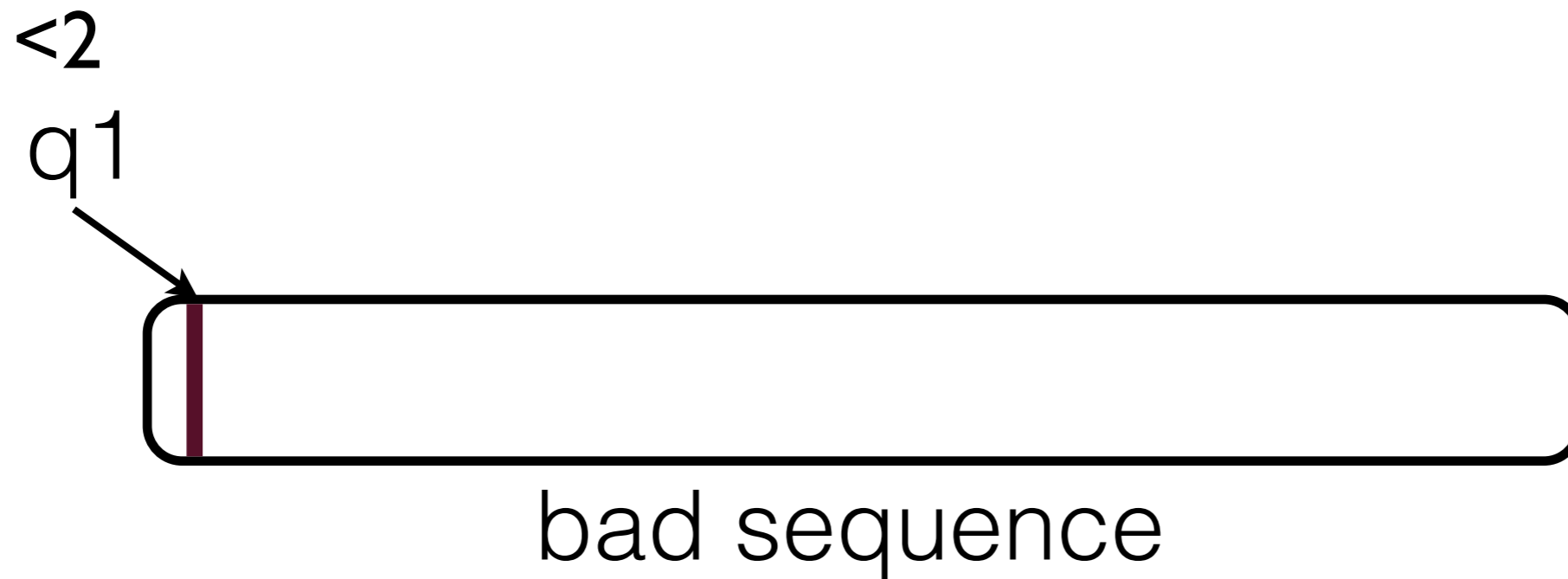
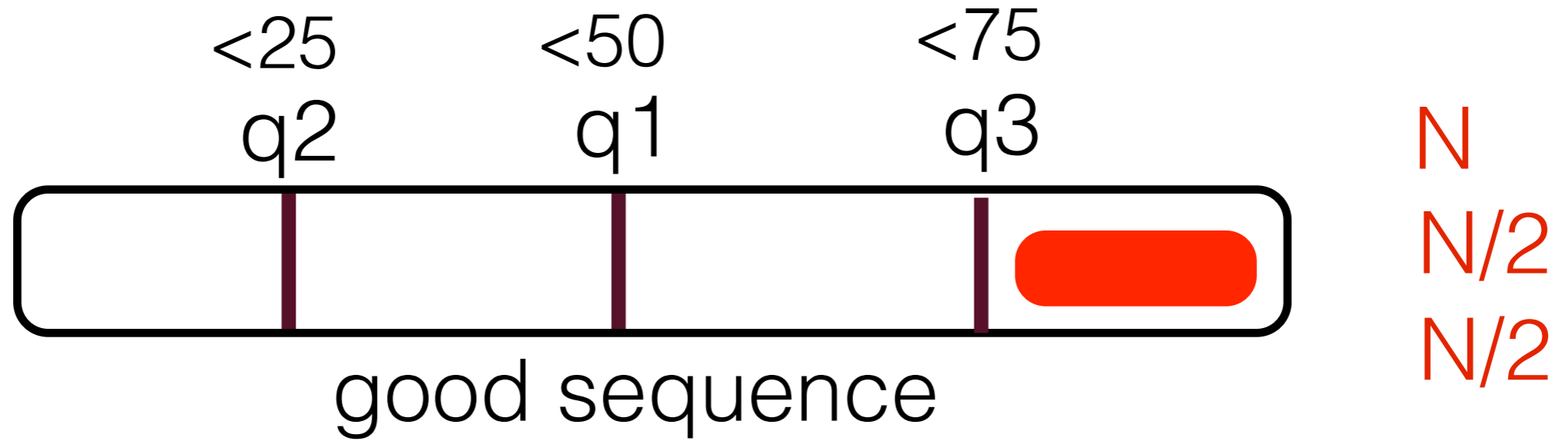
column with 100 unique integers [1,100]



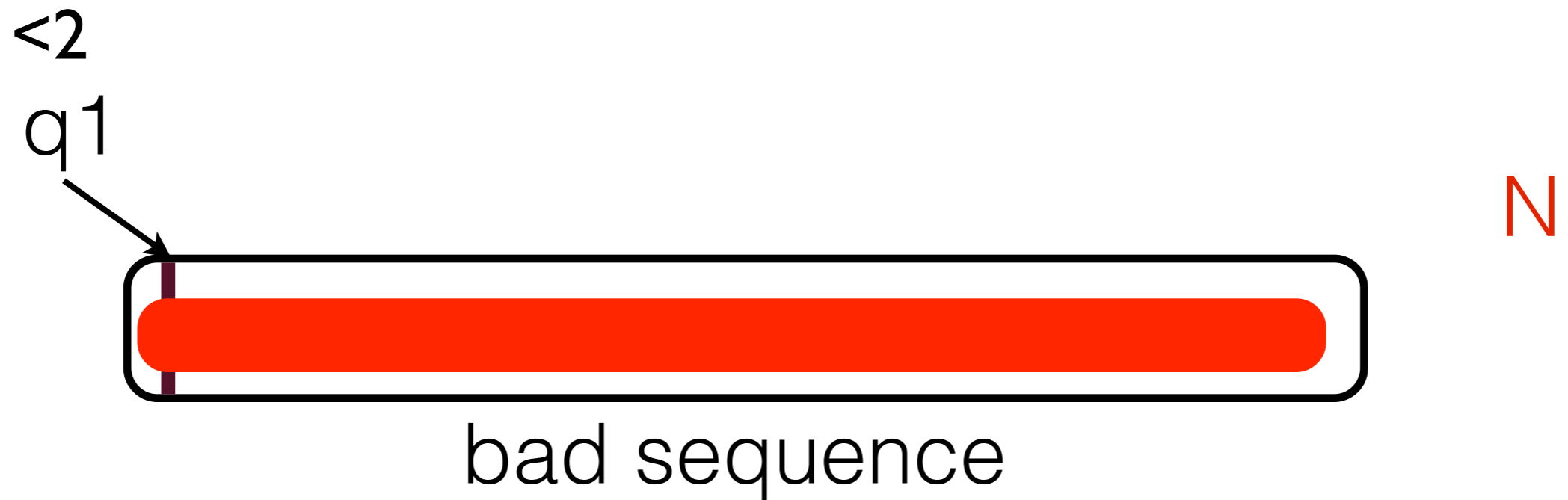
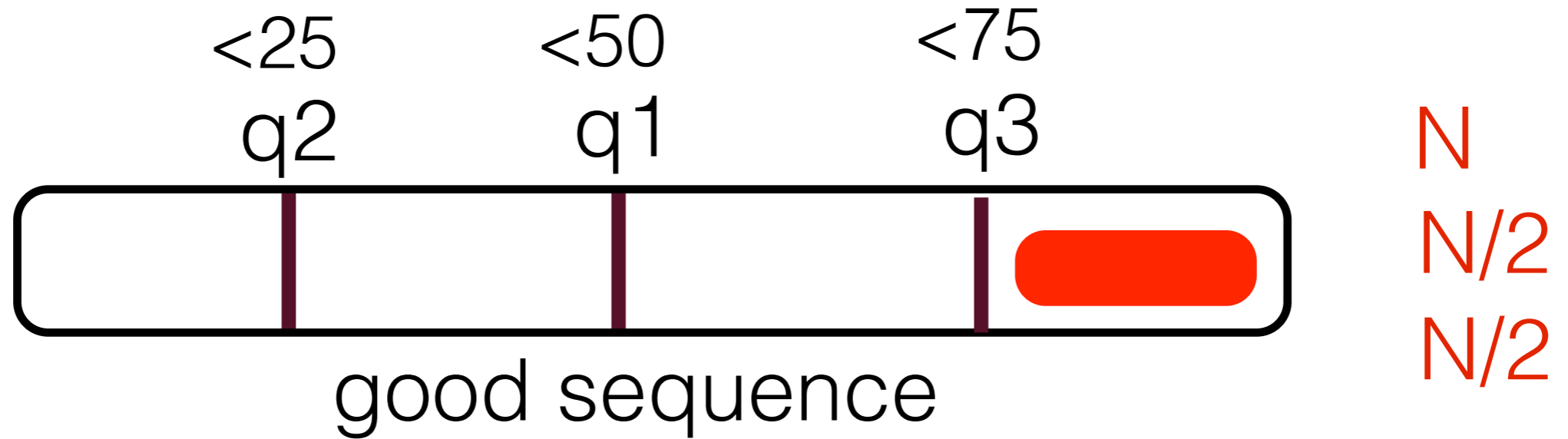
column with 100 unique integers [1,100]



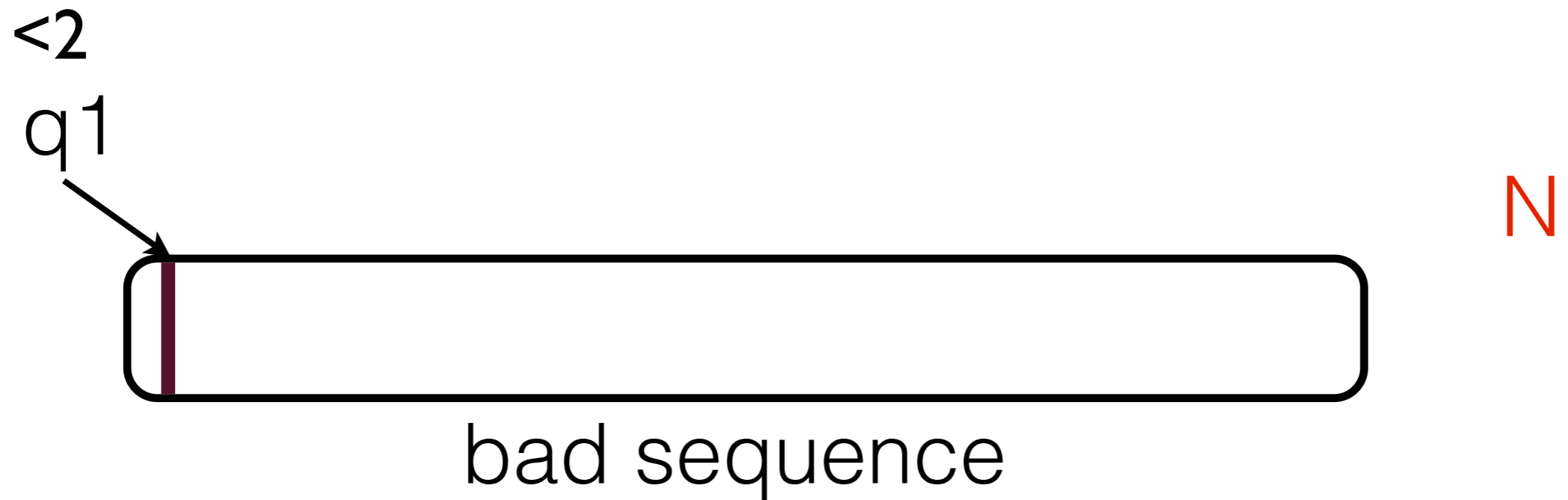
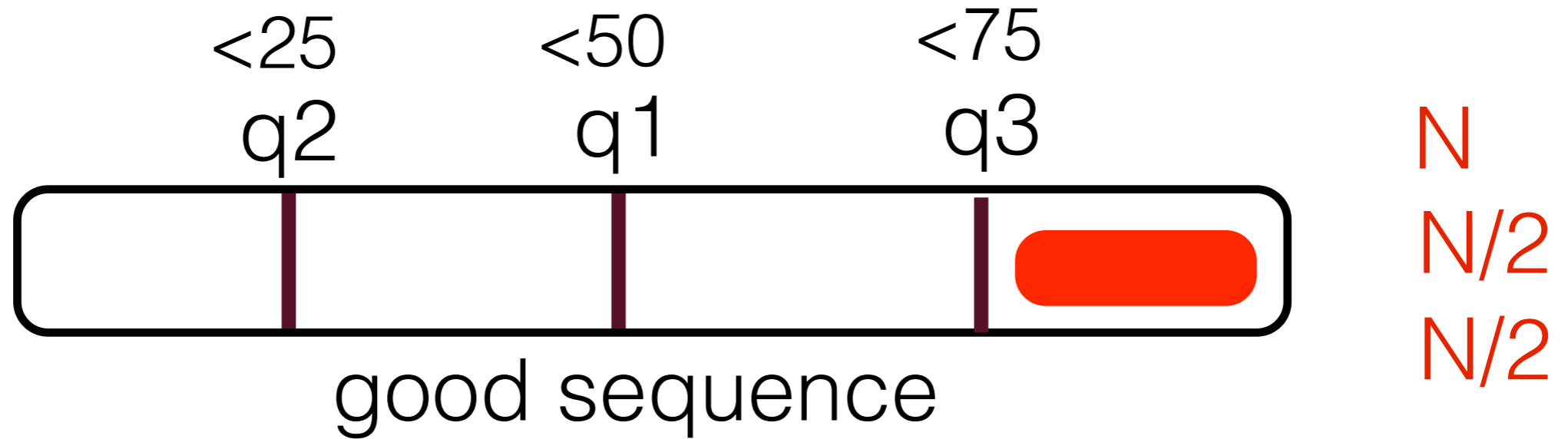
column with 100 unique integers [1,100]



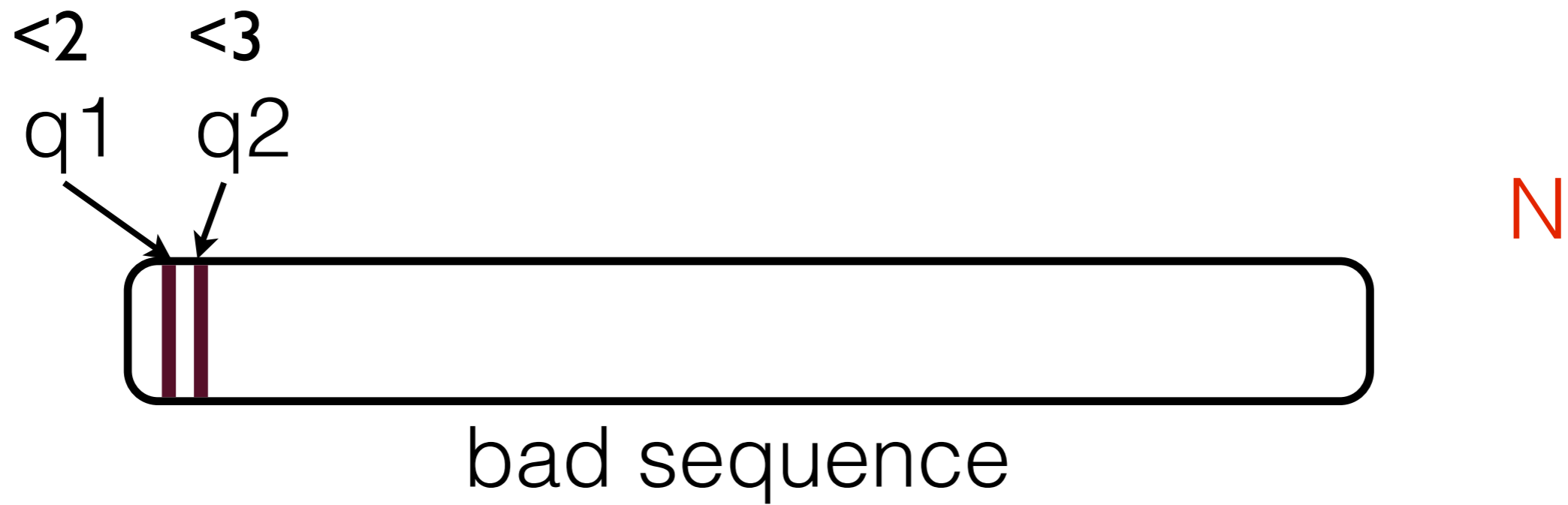
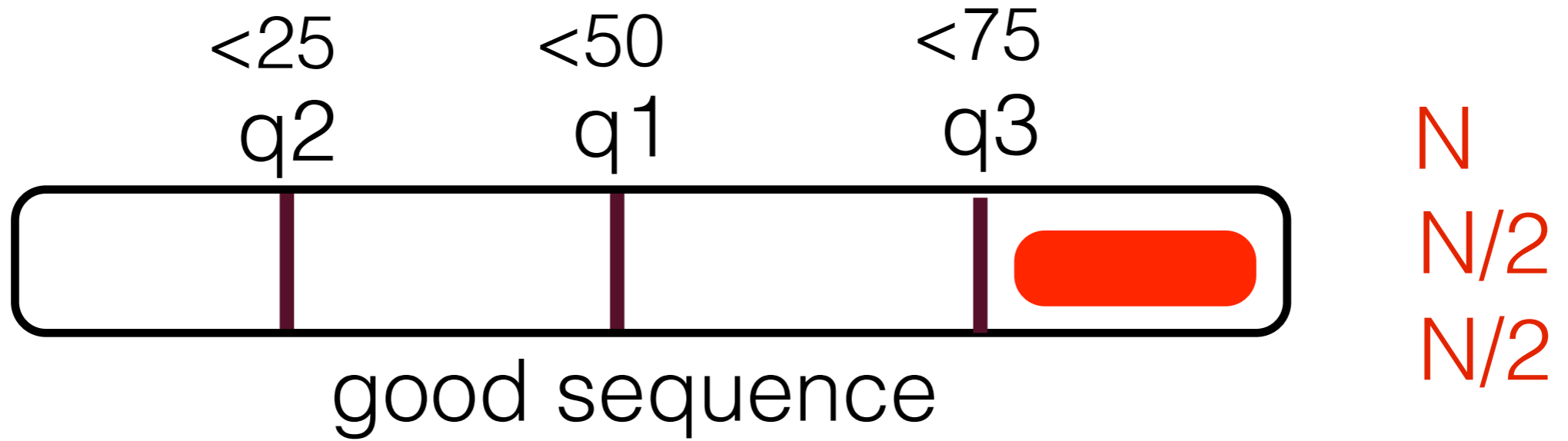
column with 100 unique integers [1,100]



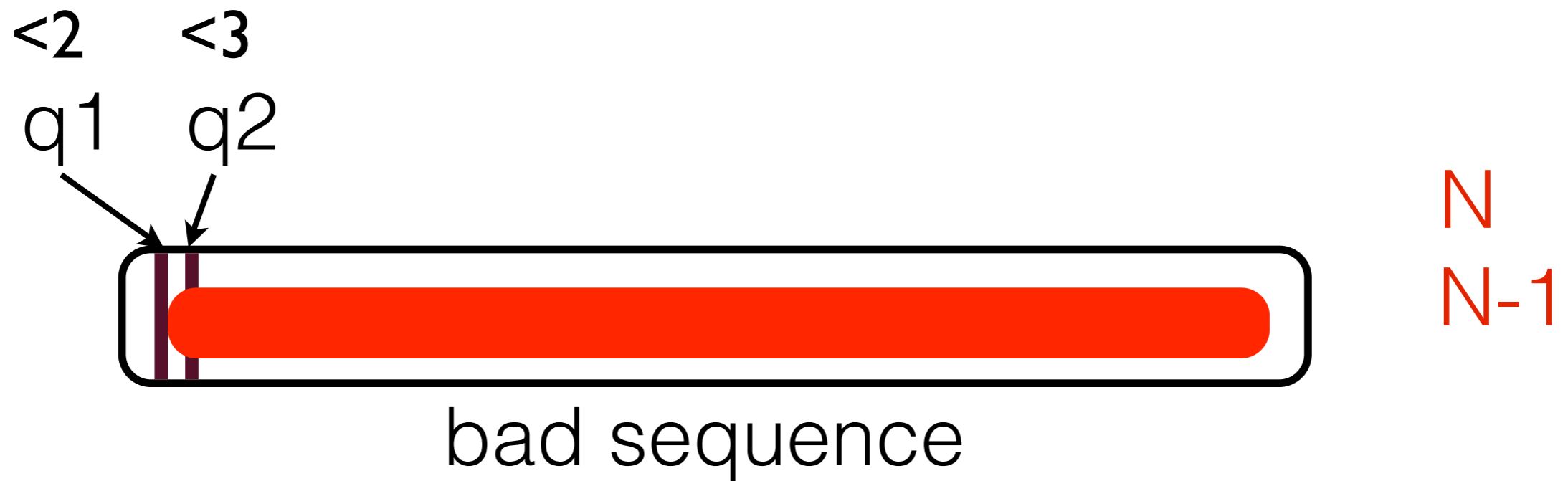
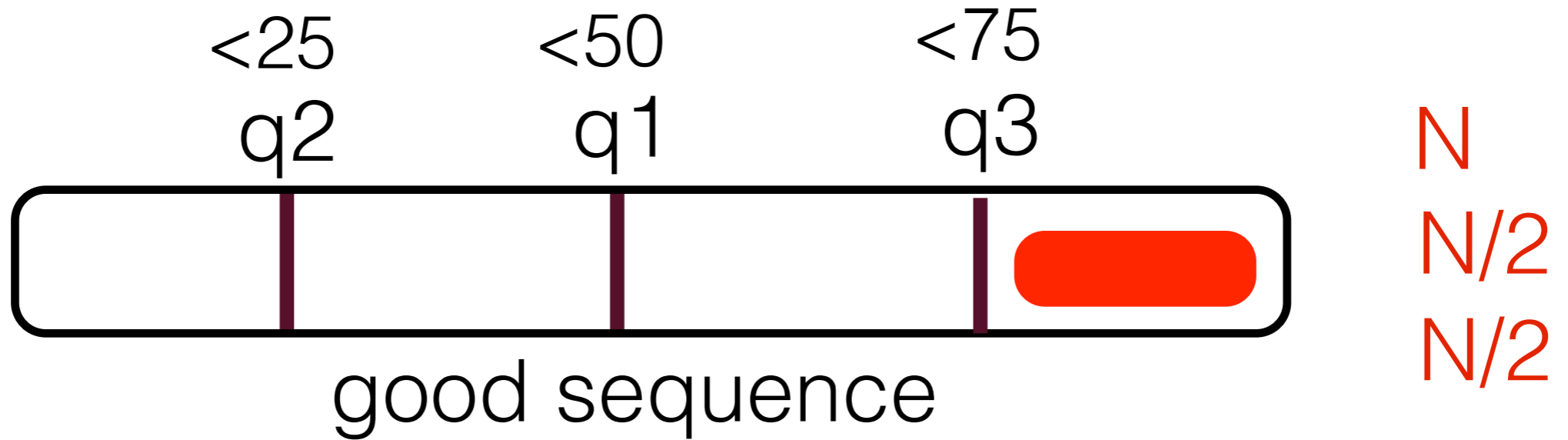
column with 100 unique integers [1,100]



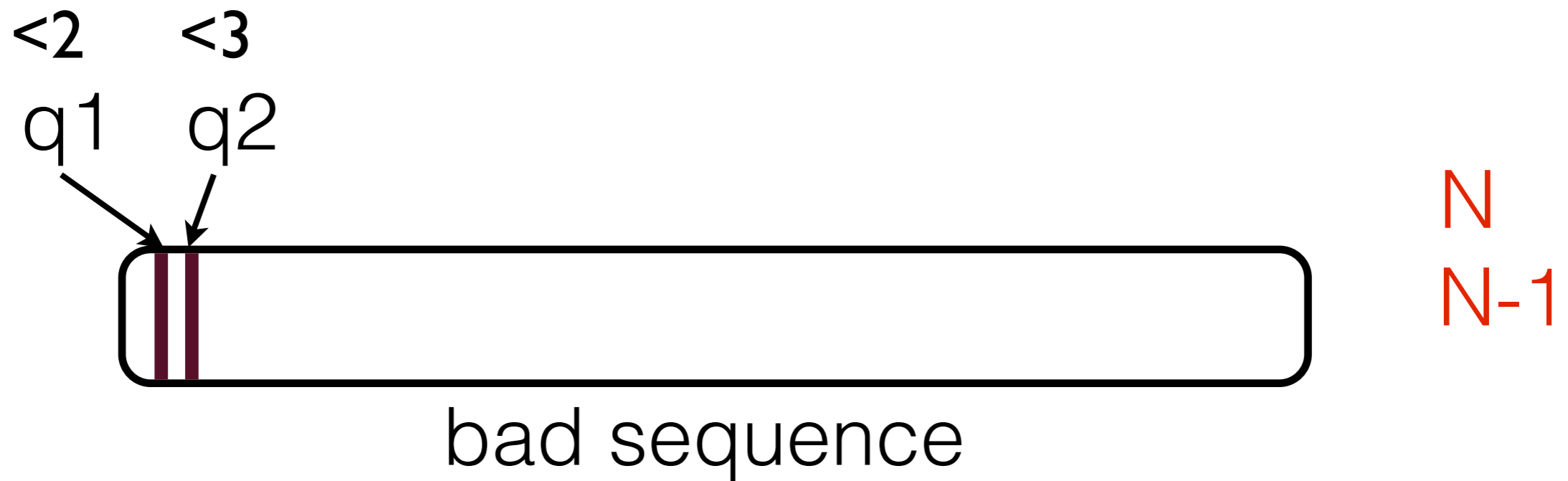
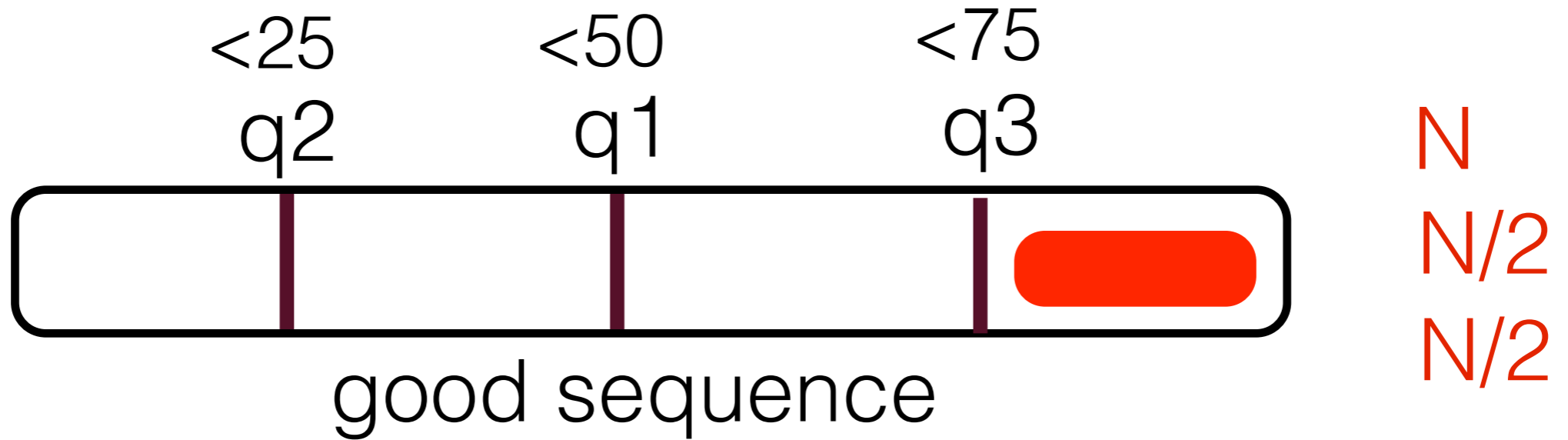
column with 100 unique integers [1,100]



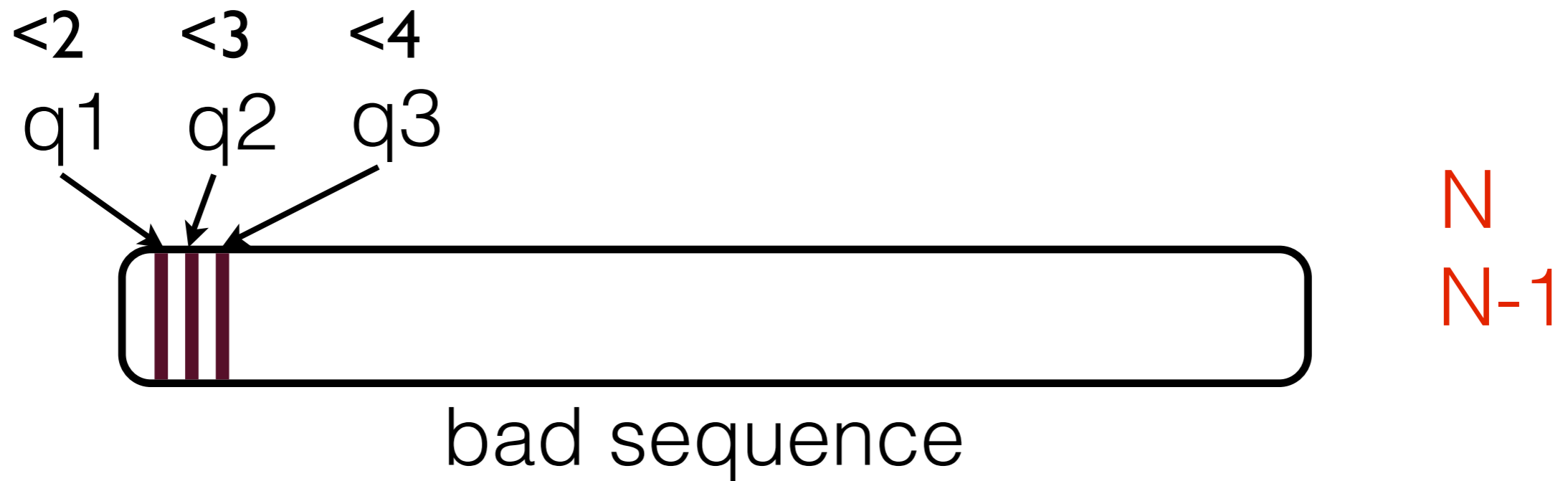
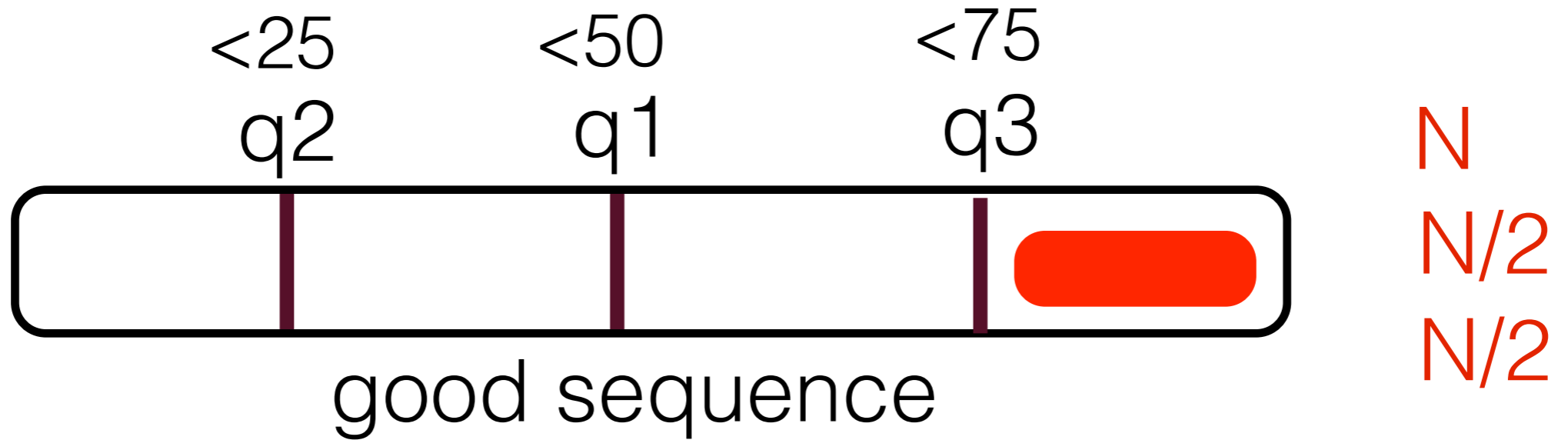
column with 100 unique integers [1,100]



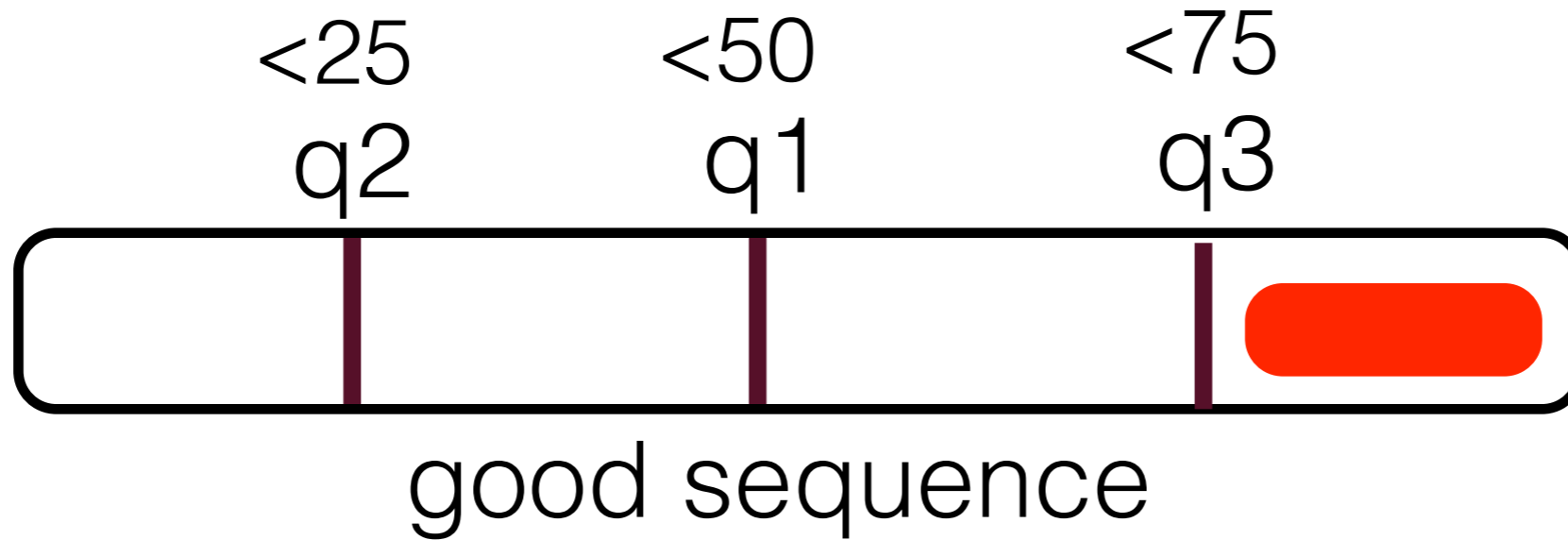
column with 100 unique integers [1,100]



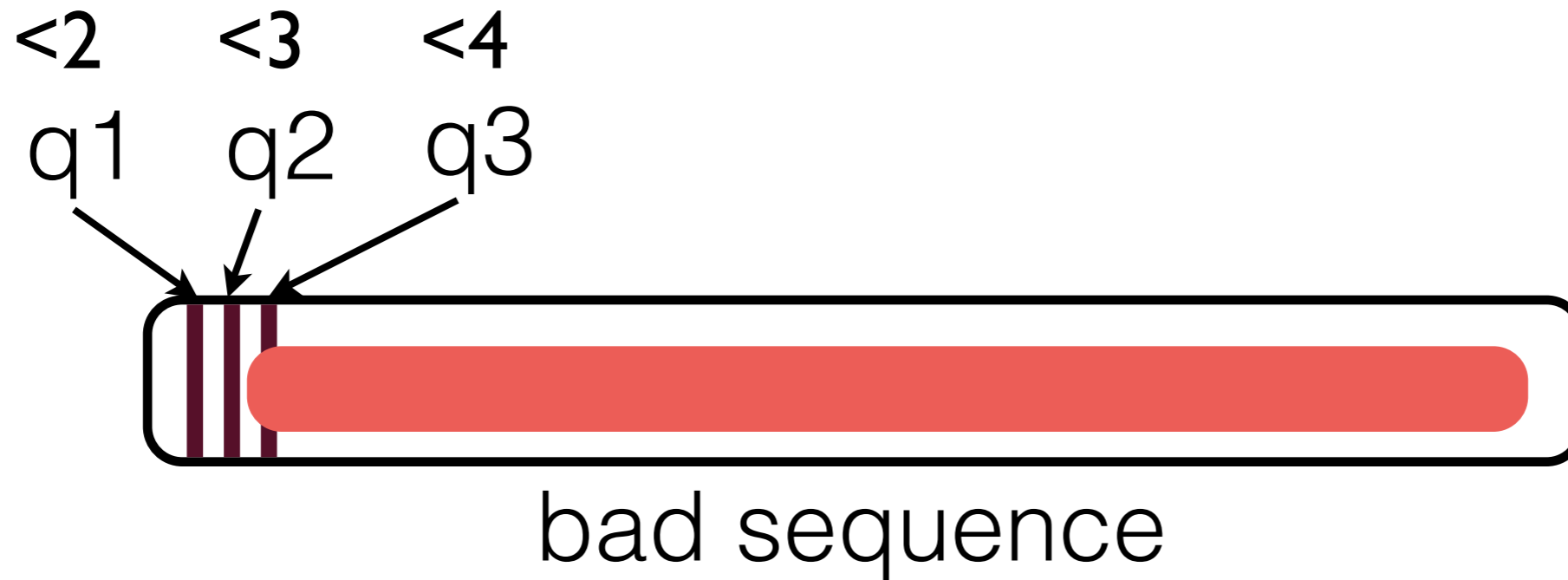
column with 100 unique integers [1,100]



column with 100 unique integers [1,100]

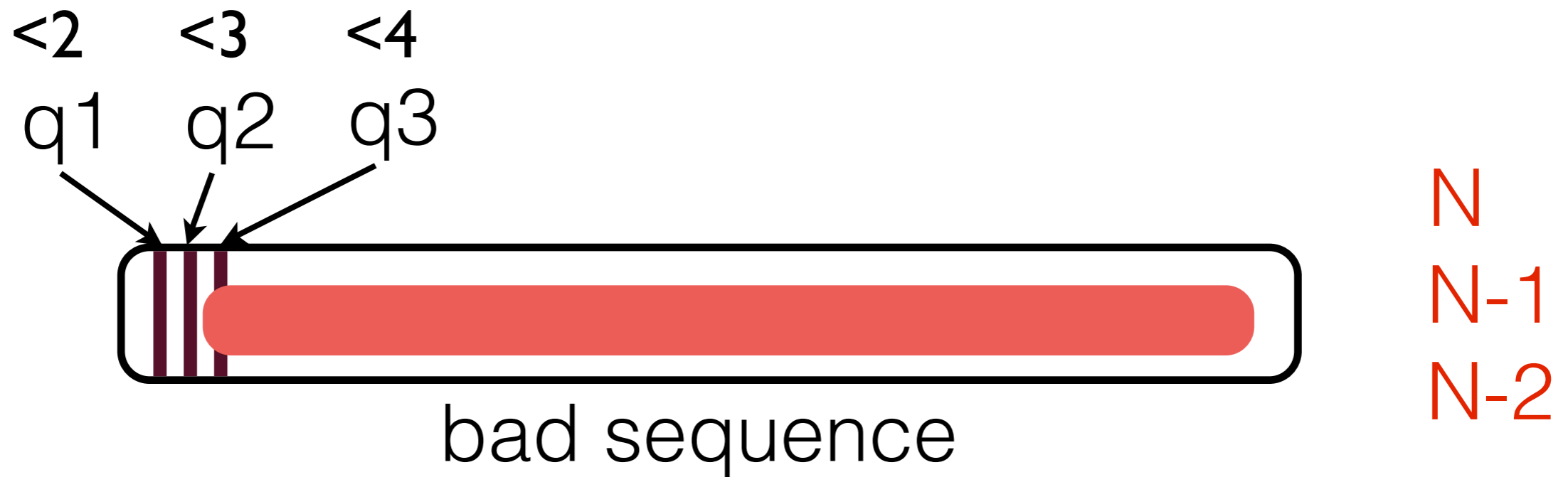


N
N/2
N/2



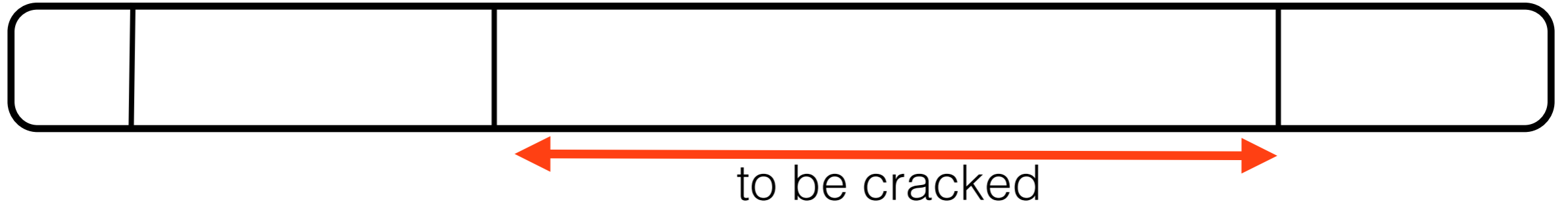
N
N-1
N-2

blindly adapting to queries is
not always a good idea

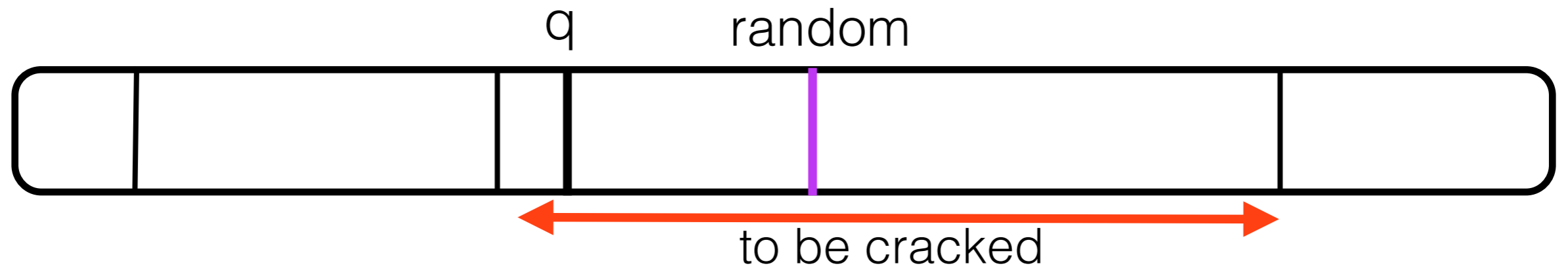


query driven

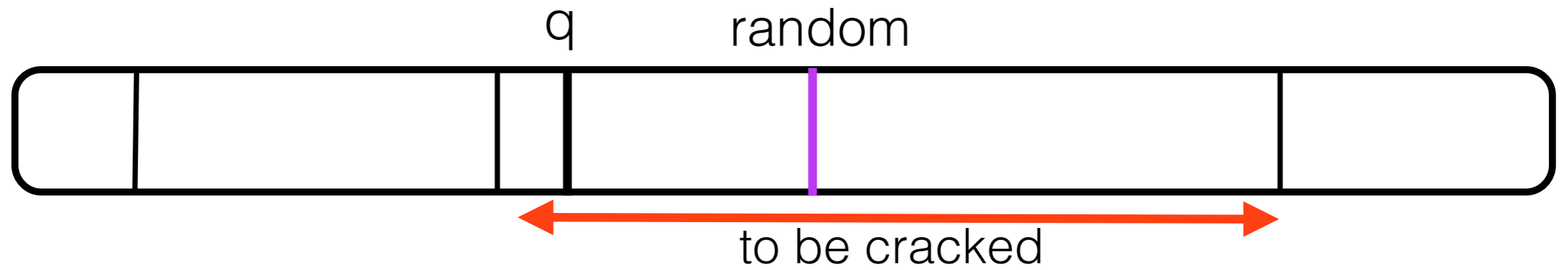
query driven



query driven

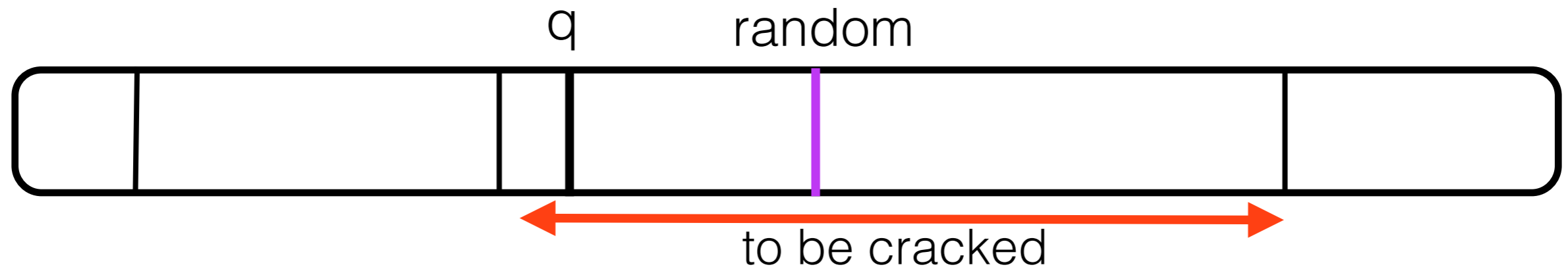


query driven



progressive cracking

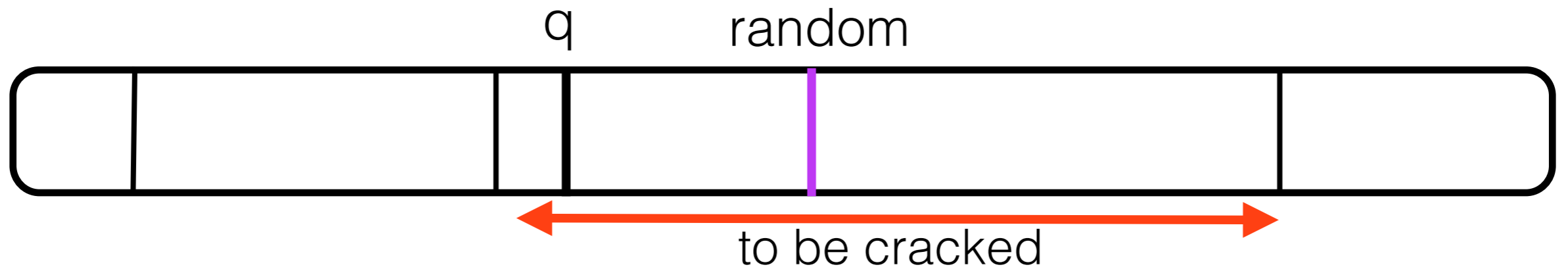
query driven



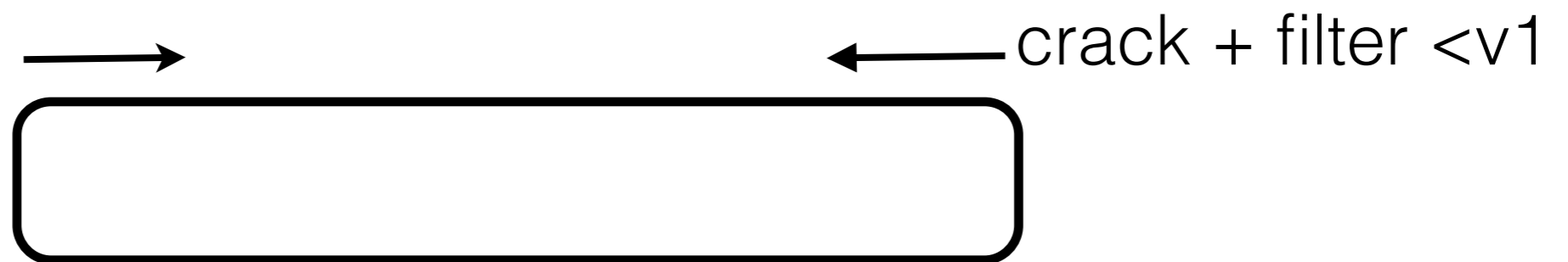
progressive cracking
q1: <v1



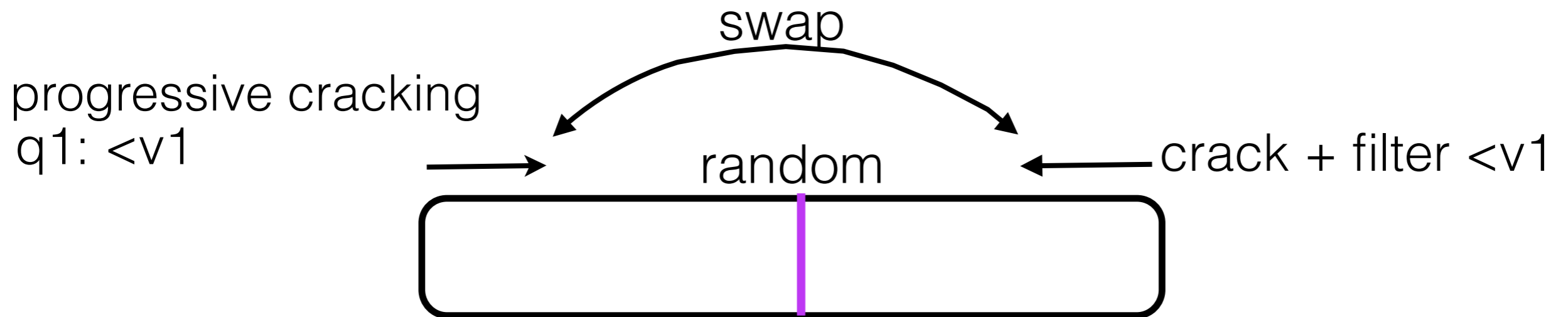
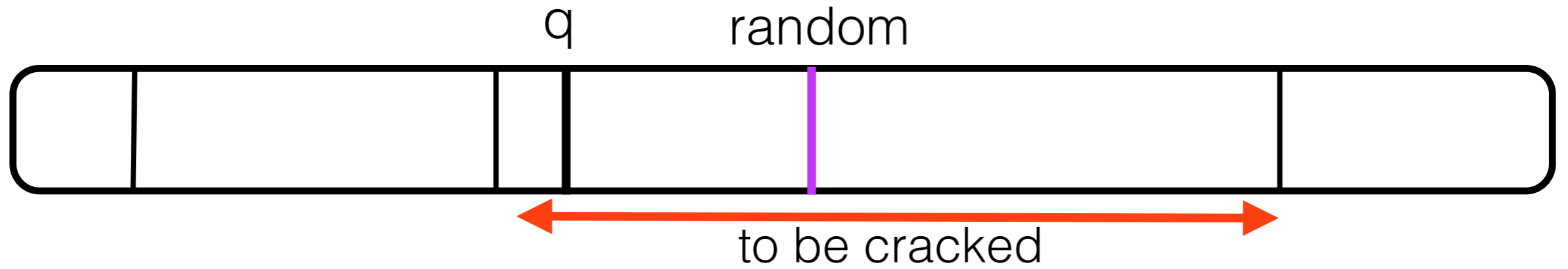
query driven



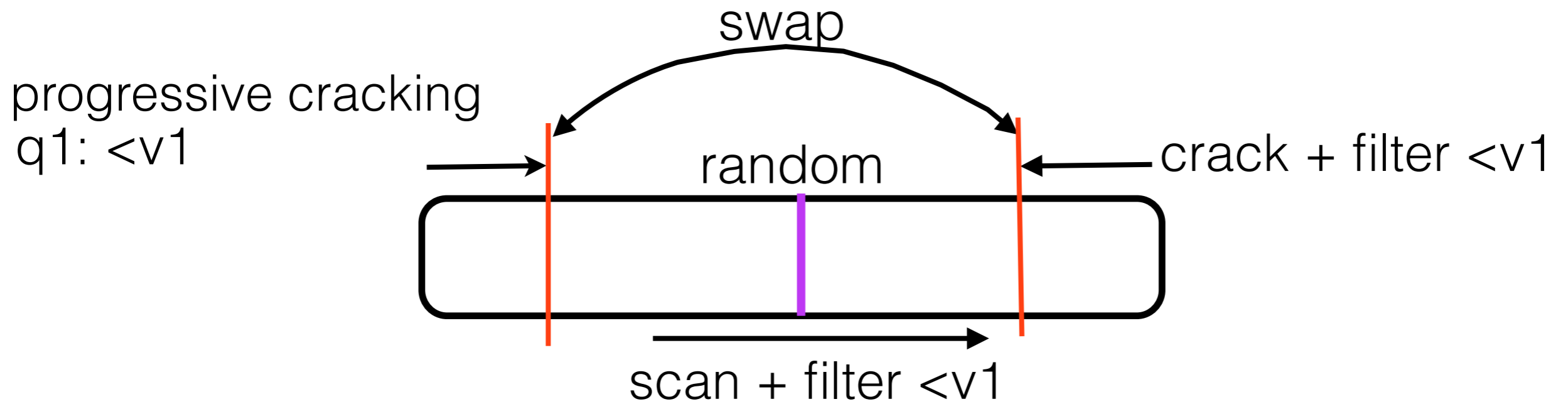
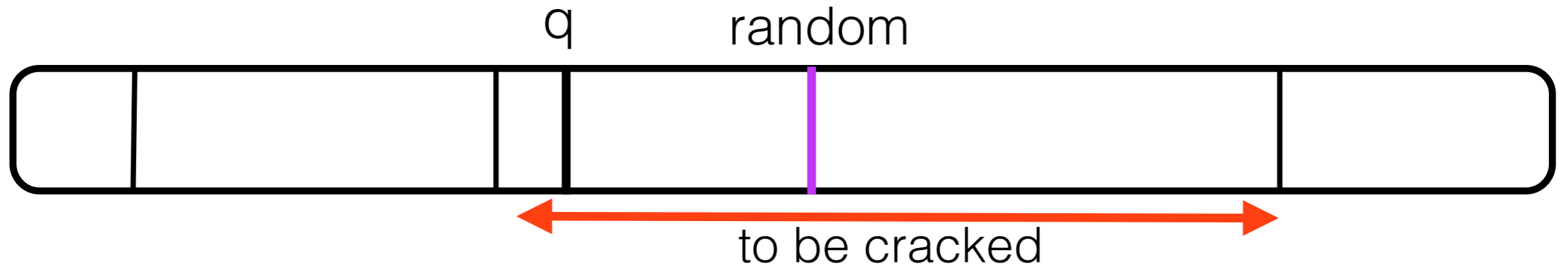
progressive cracking
q1: <v1



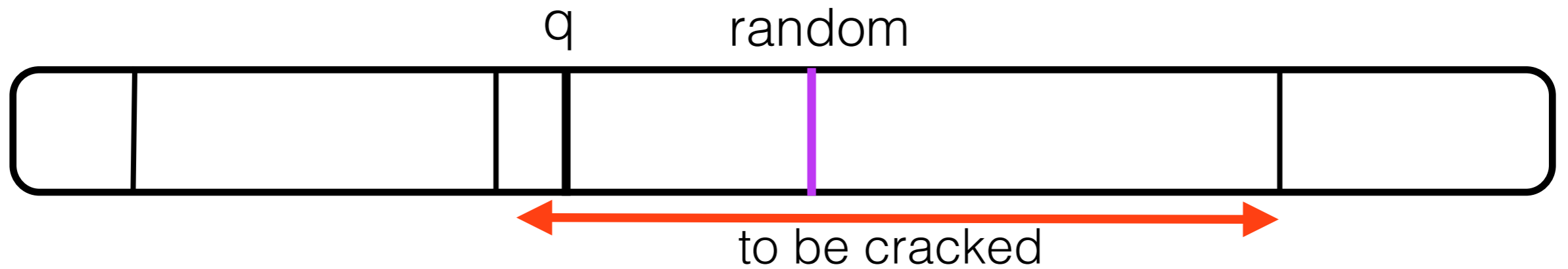
query driven



query driven



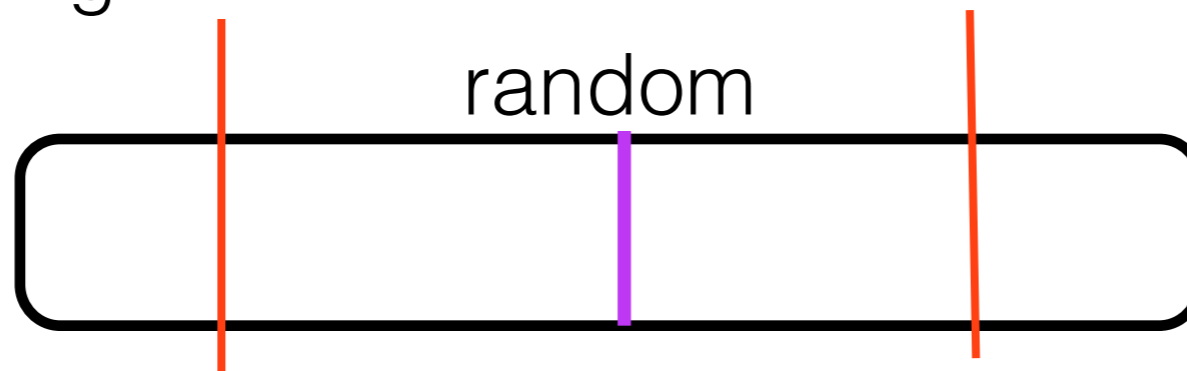
query driven



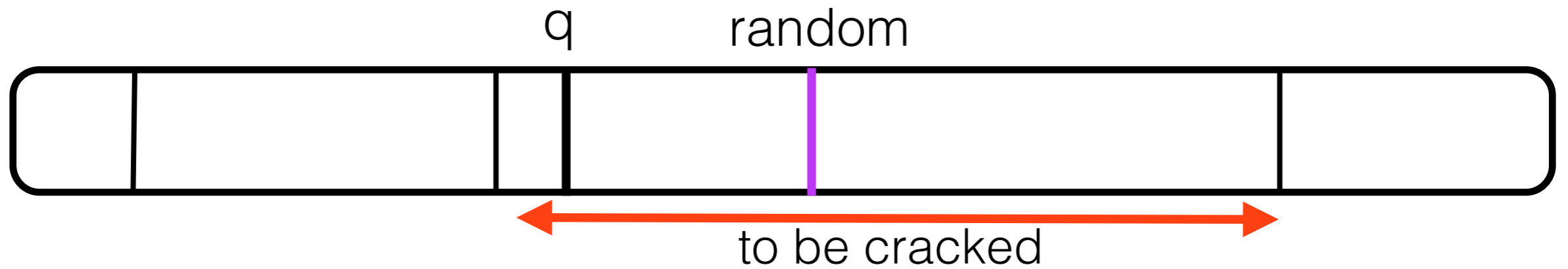
progressive cracking

q1: <v1

q2: <v2



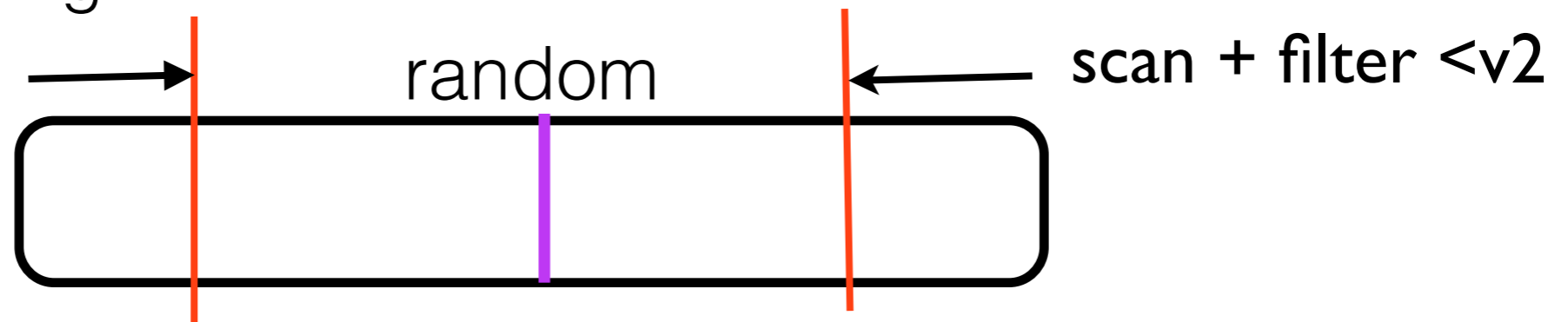
query driven



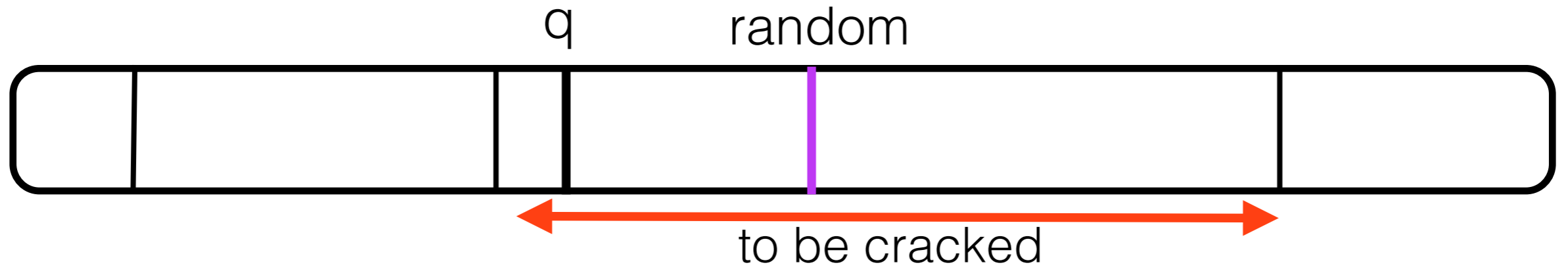
progressive cracking

q1: $<v1$

q2: $<v2$



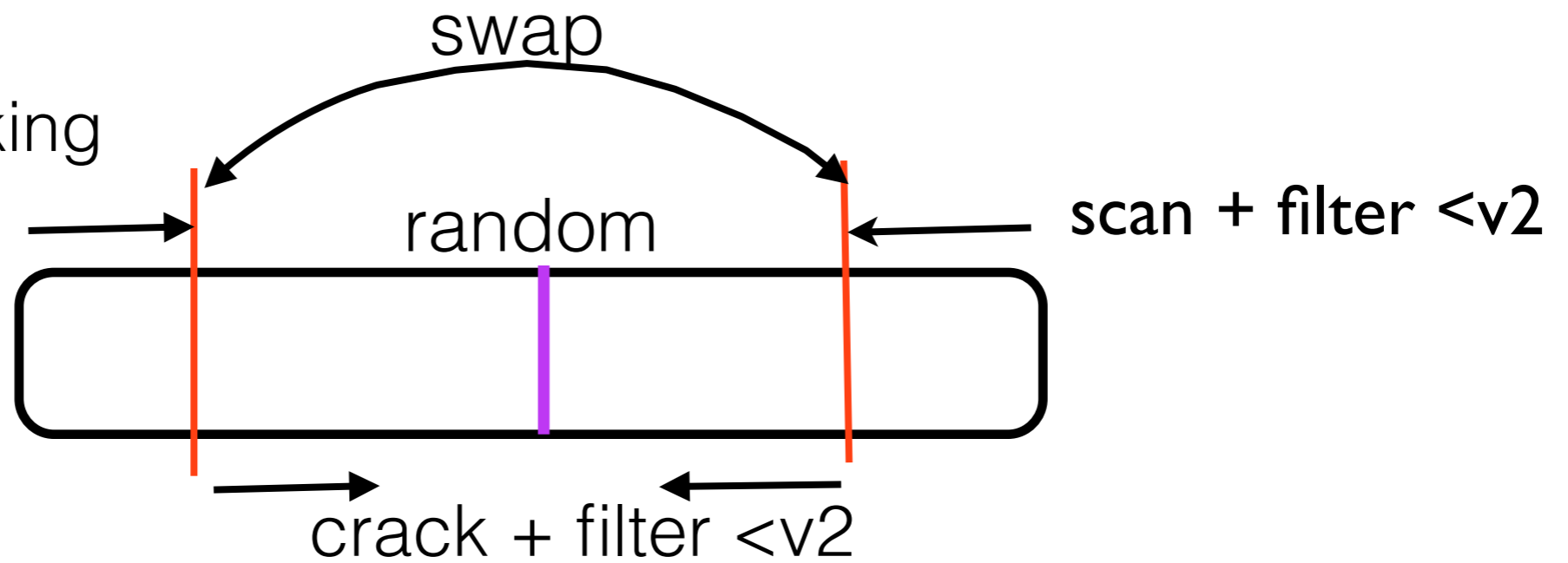
query driven



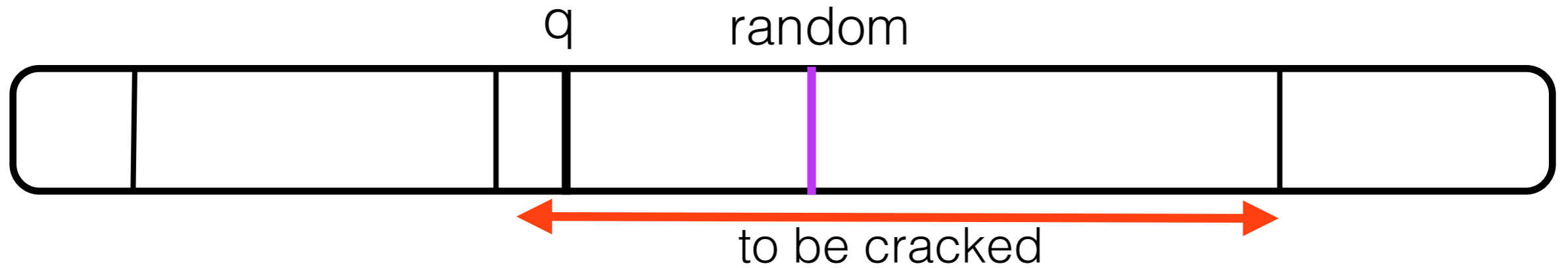
progressive cracking

q1: $<v1$

q2: $<v2$



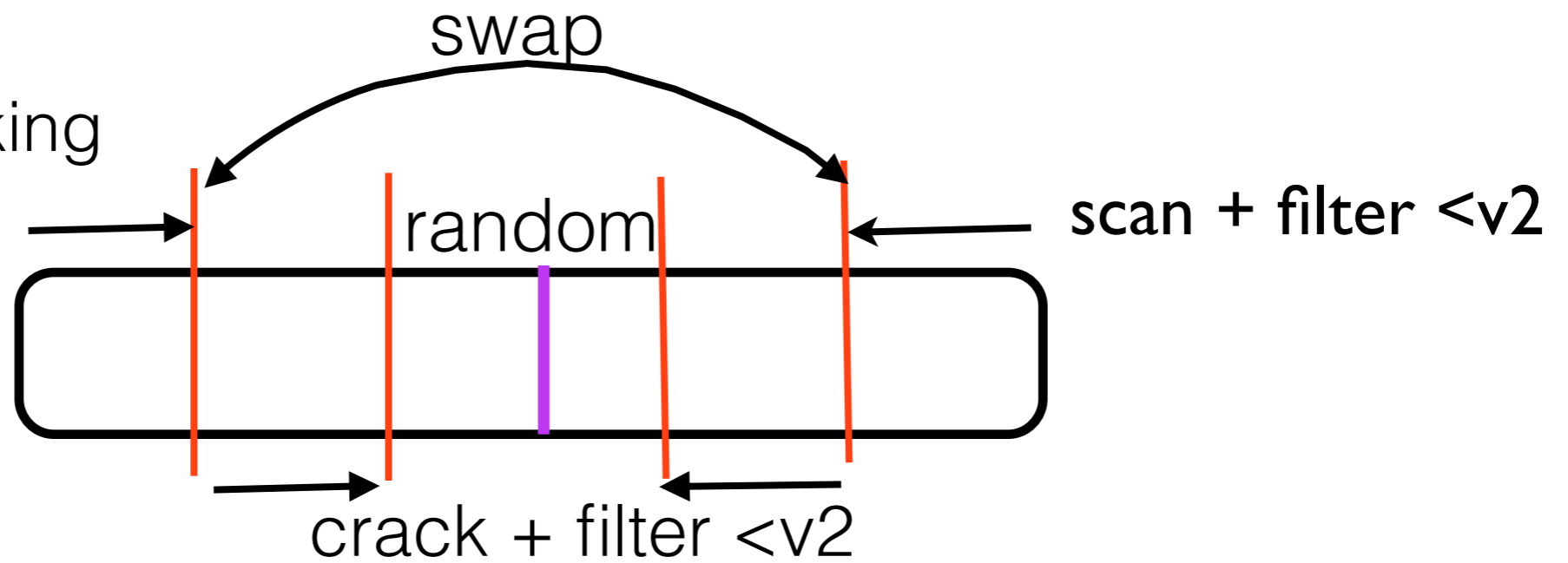
query driven



progressive cracking

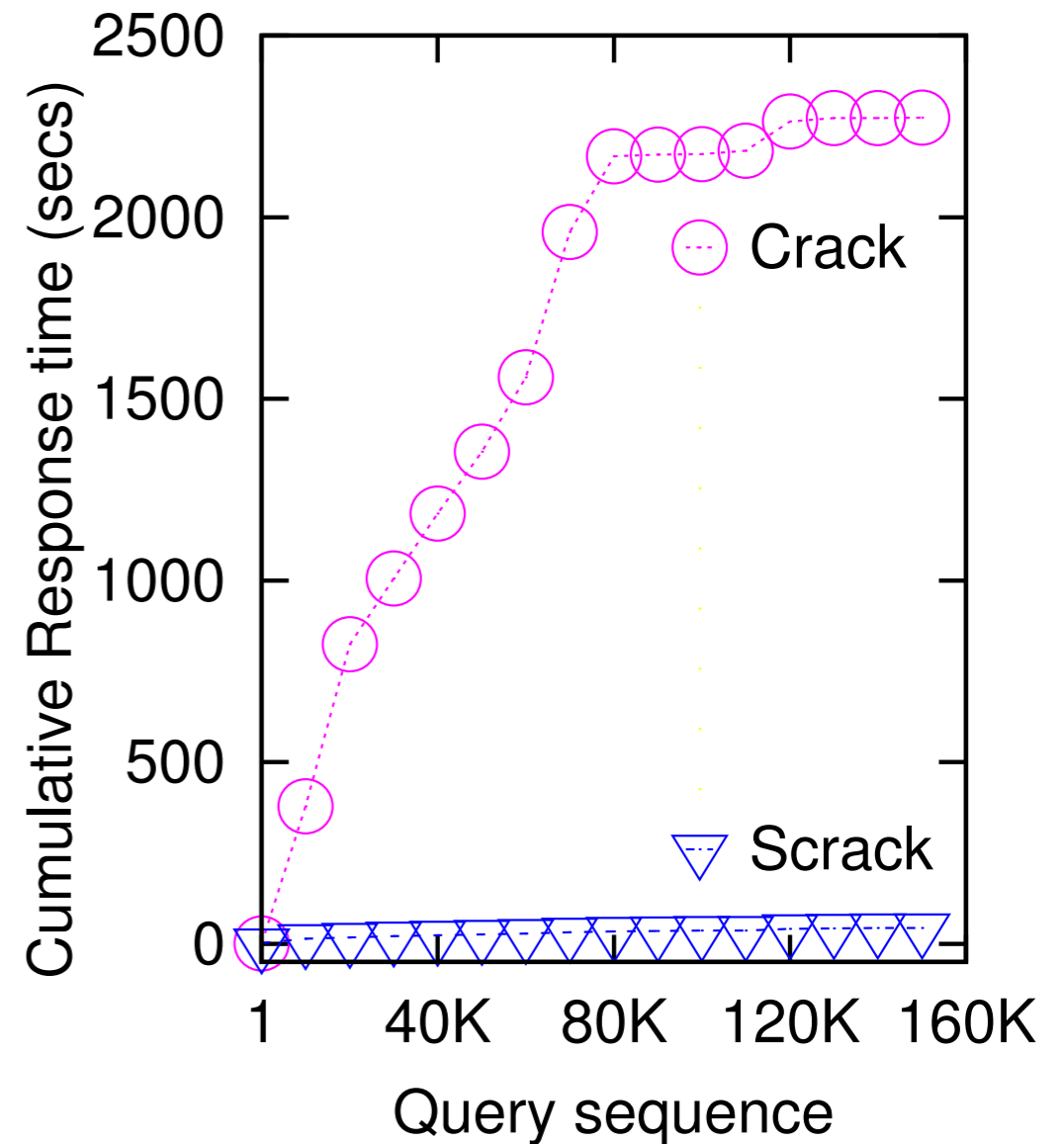
q1: $<v1$

q2: $<v2$



cracking on Skyserver (4TB)

(Sloan Digital Sky Survey, www.sdss.org)



cracking answers 160.000 queries
while full indexing is still half way creating one index

concurrency control

problem: read queries become write queries!

goal: be able to crack for multiple queries in parallel

traditional indexing

adaptive indexing

write queries

traditional indexing

read queries

adaptive indexing

change index contents
and structure

write queries

traditional indexing

only index structure
changes

read queries

adaptive indexing

no need for traditional **locks** = too heavy

short term **latches** = fast and release quickly

change index contents
and structure

write queries

traditional indexing

only index structure
changes

read queries

adaptive indexing

all or nothing

change index contents
and structure

write queries

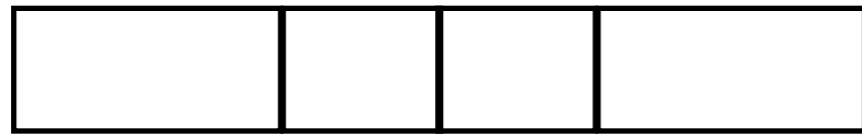
traditional indexing

incremental and optional

only index structure
changes

read queries

adaptive indexing



v1 v2 v3

all or nothing

change index contents
and structure

write queries

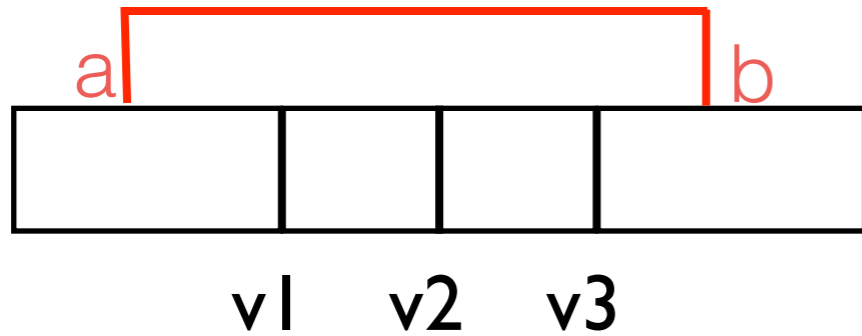
traditional indexing

incremental and optional

only index structure
changes

read queries

adaptive indexing



all or nothing

change index contents
and structure

write queries

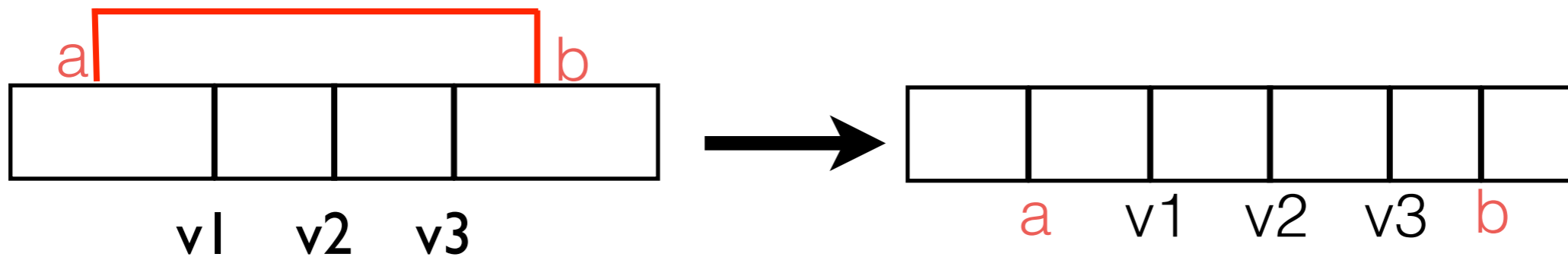
traditional indexing

incremental and optional

only index structure
changes

read queries

adaptive indexing



all or nothing

change index contents
and structure

write queries

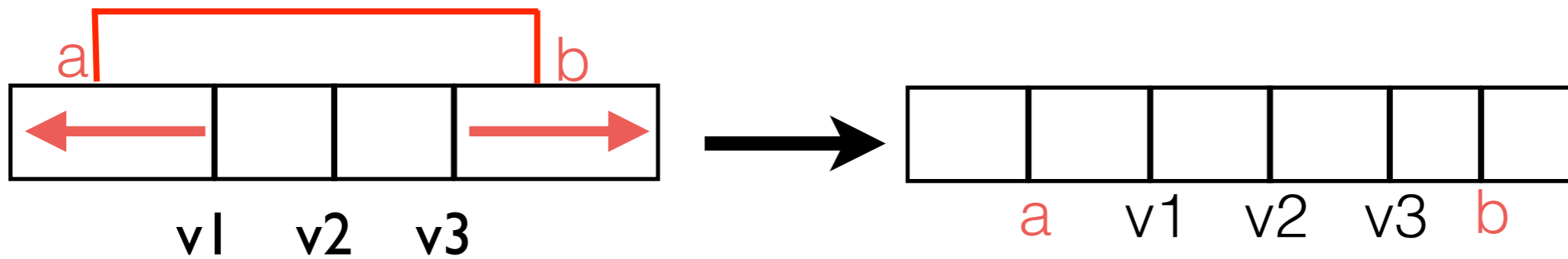
traditional indexing

incremental and optional

only index structure
changes

read queries

adaptive indexing



all or nothing

change index contents
and structure

write queries

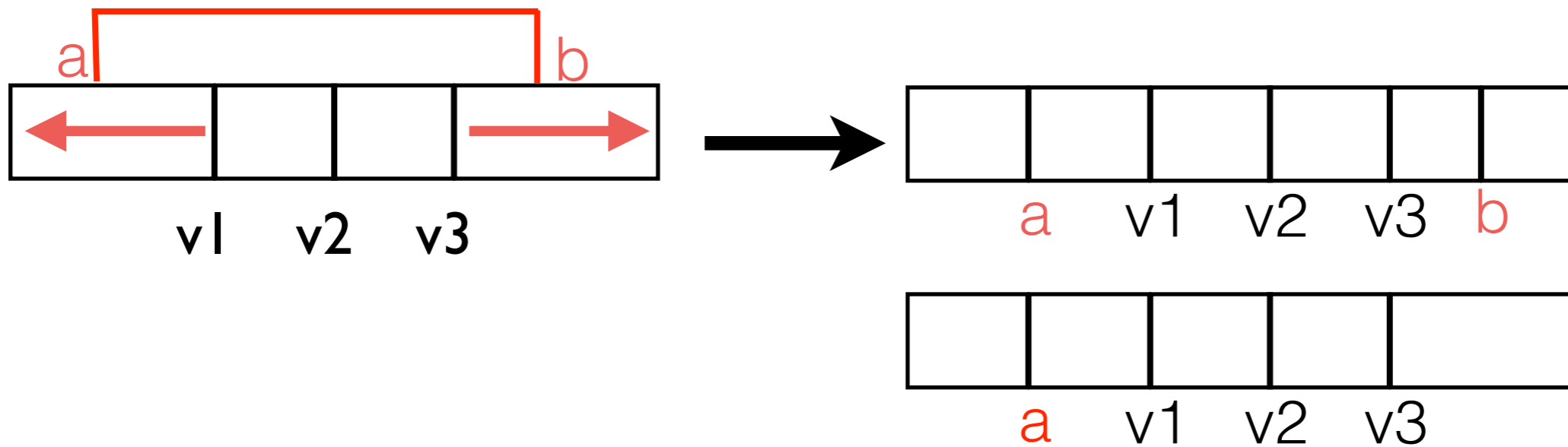
traditional indexing

incremental and optional

only index structure
changes

read queries

adaptive indexing



all or nothing

change index contents
and structure

write queries

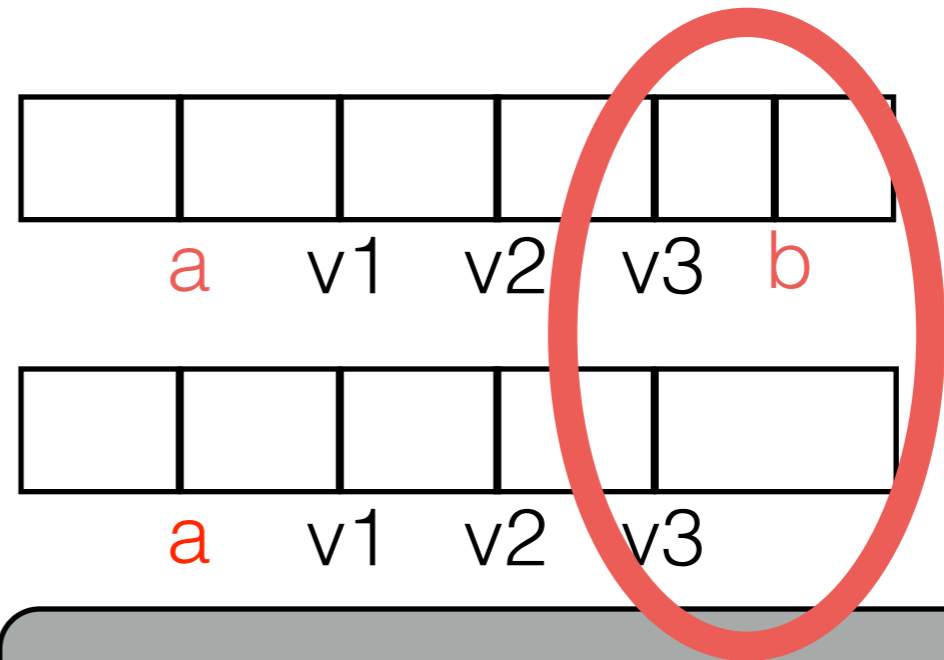
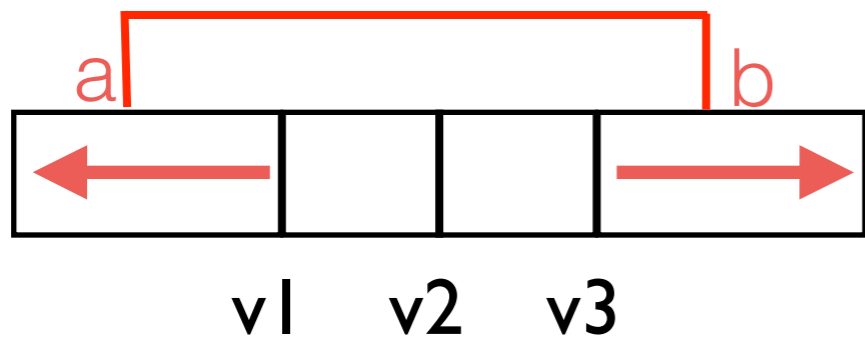
traditional indexing

incremental and optional

only index structure
changes

read queries

adaptive indexing



all or nothing

change index contents and structure

write queries

traditional indexing

incremental and optional

only index structure changes

read queries

adaptive indexing

impact stable storage

all or nothing

change index contents
and structure

write queries

traditional indexing

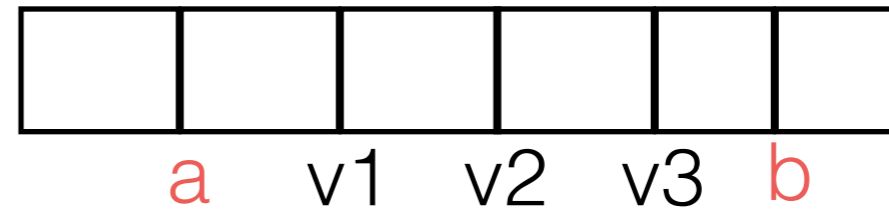
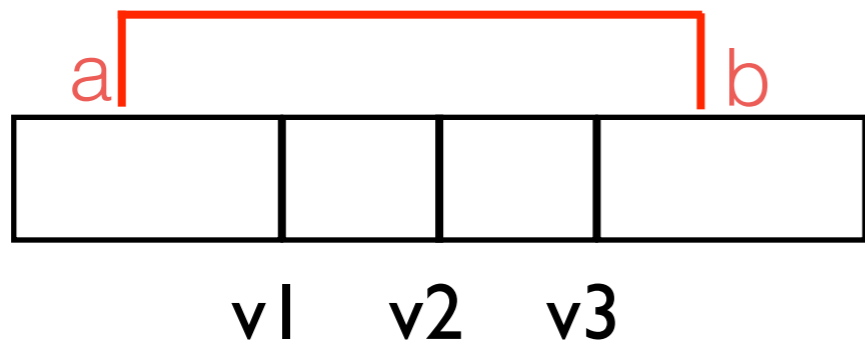
stable storage optional

incremental and optional

only index structure
changes

read queries

adaptive indexing



impact stable storage

all or nothing

change index contents and structure

write queries

traditional indexing

stable storage optional

incremental and optional

only index structure changes

read queries

adaptive indexing

need to serialize

impact stable storage

all or nothing

change index contents
and structure

write queries

traditional indexing

can execute in any order

stable storage optional

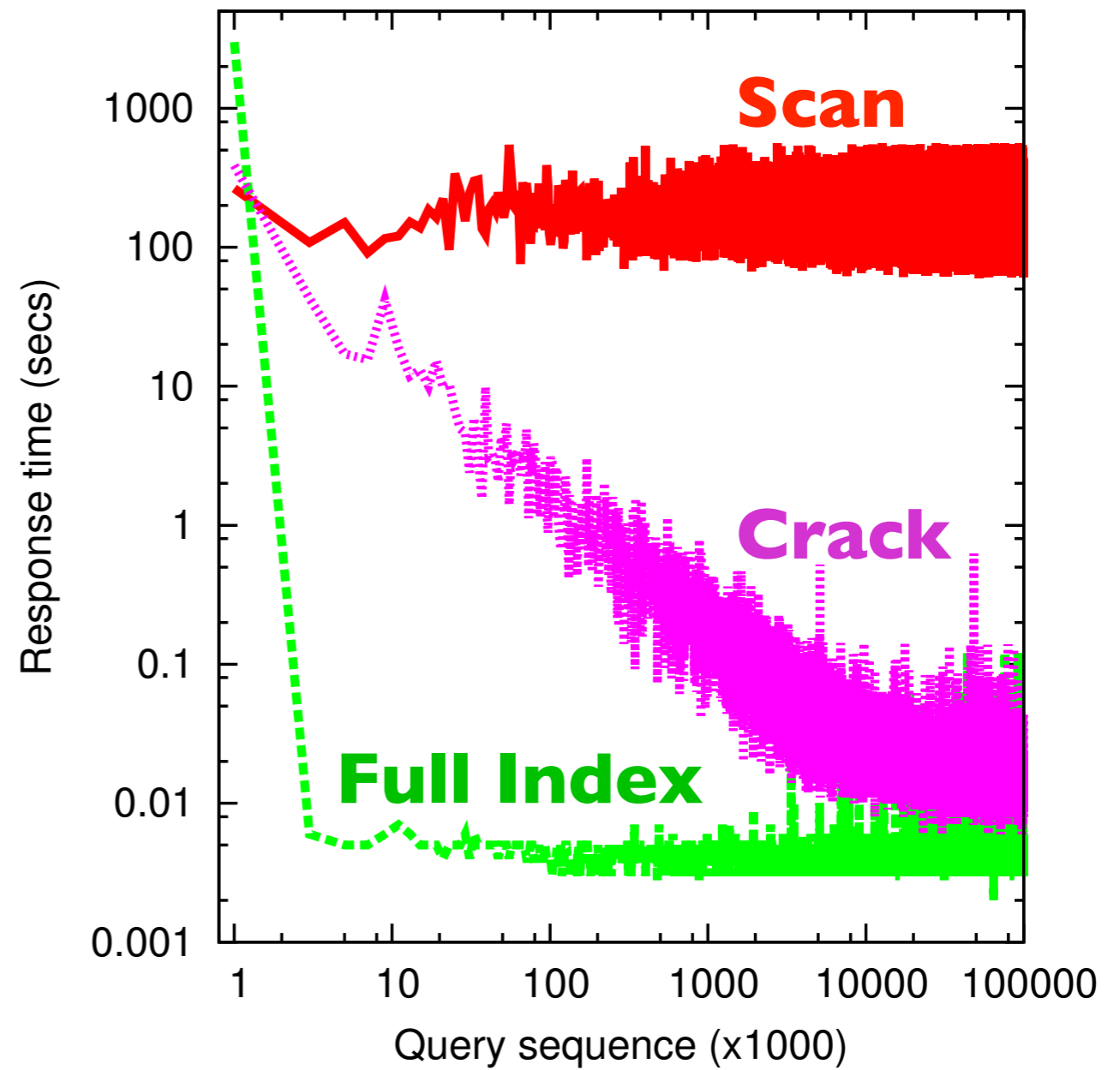
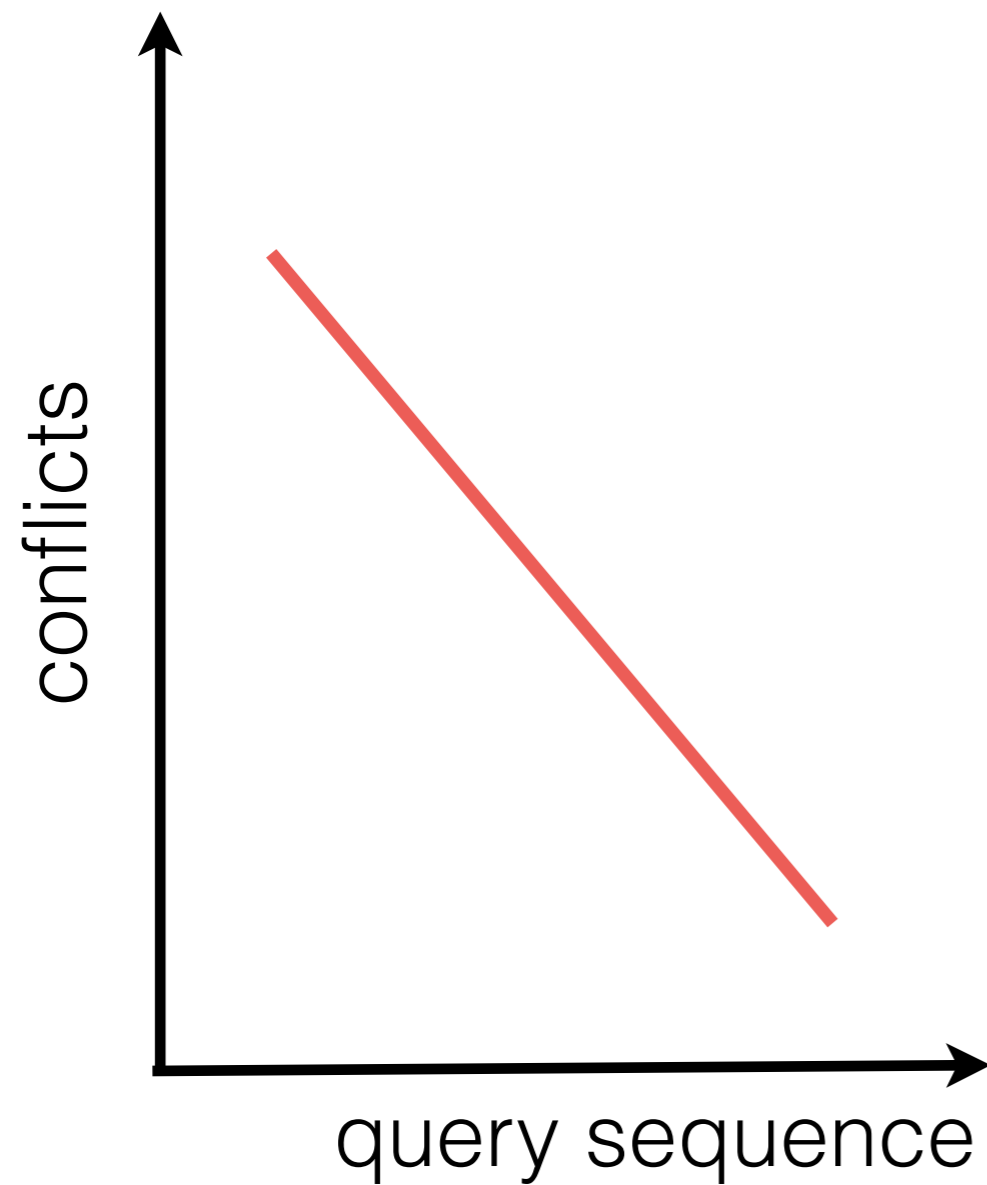
incremental and optional

only index structure
changes

read queries

adaptive indexing

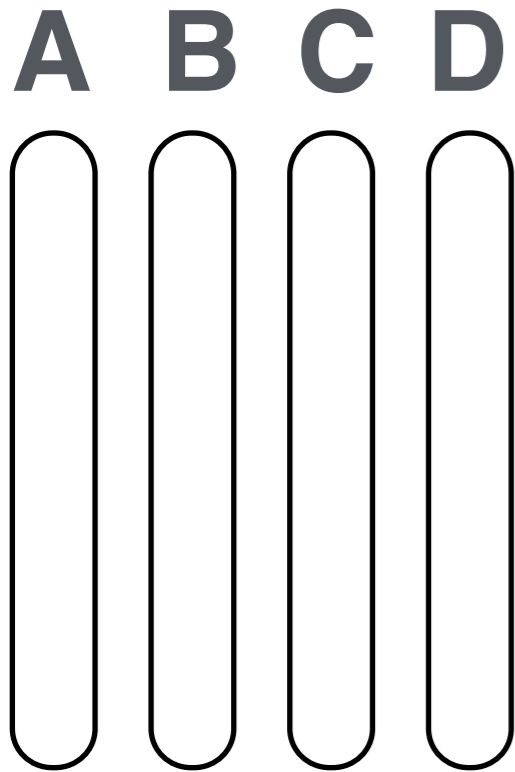
fewer conflicts as we adapt



select min(C) from R where $A < 10$ & $B < 20$



disk

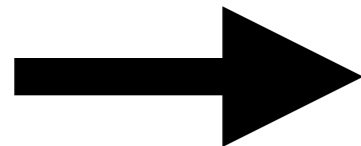
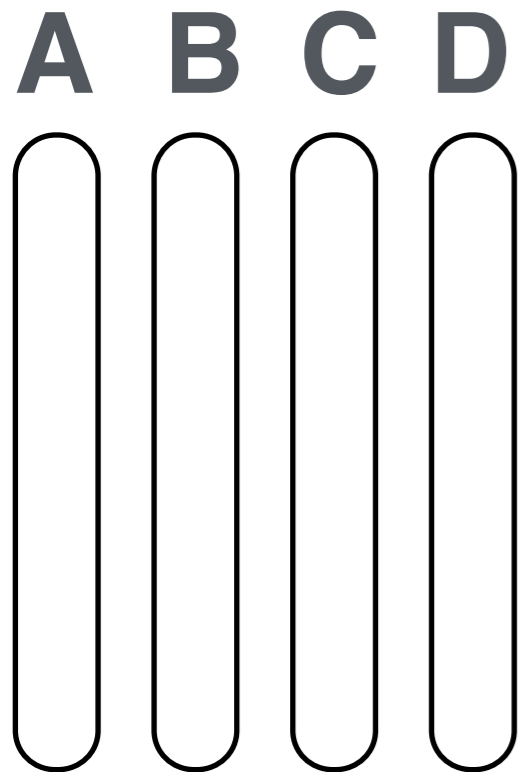


select $\min(C)$ from R where $A < 10$ & $B < 20$



disk

memory



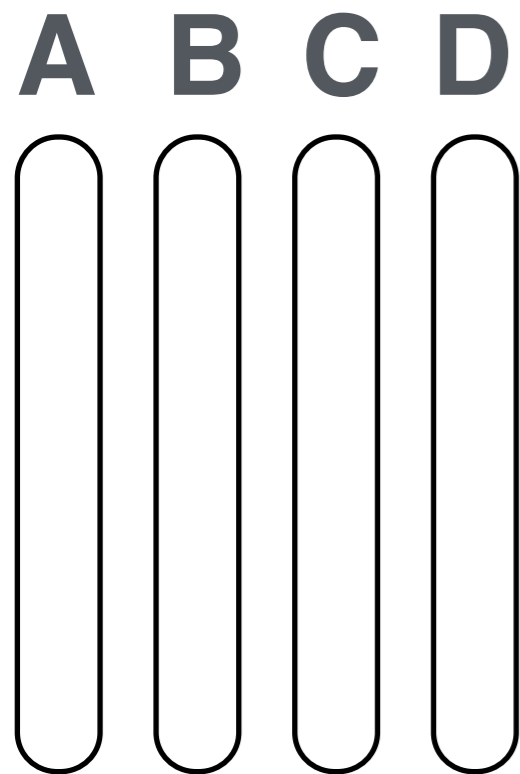
$A < 10$



select $\min(C)$ from R where $A < 10$ & $B < 20$

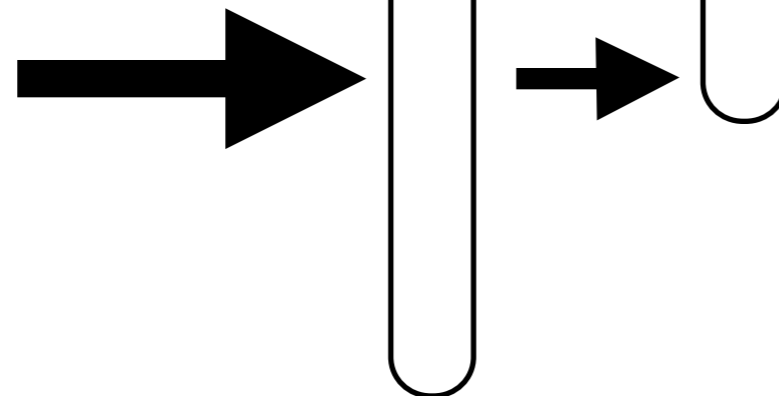


disk



memory

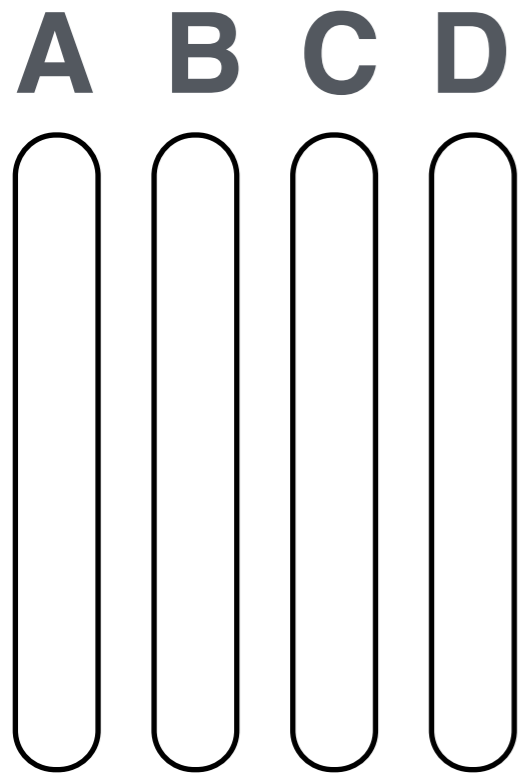
$A < 10$ IDs



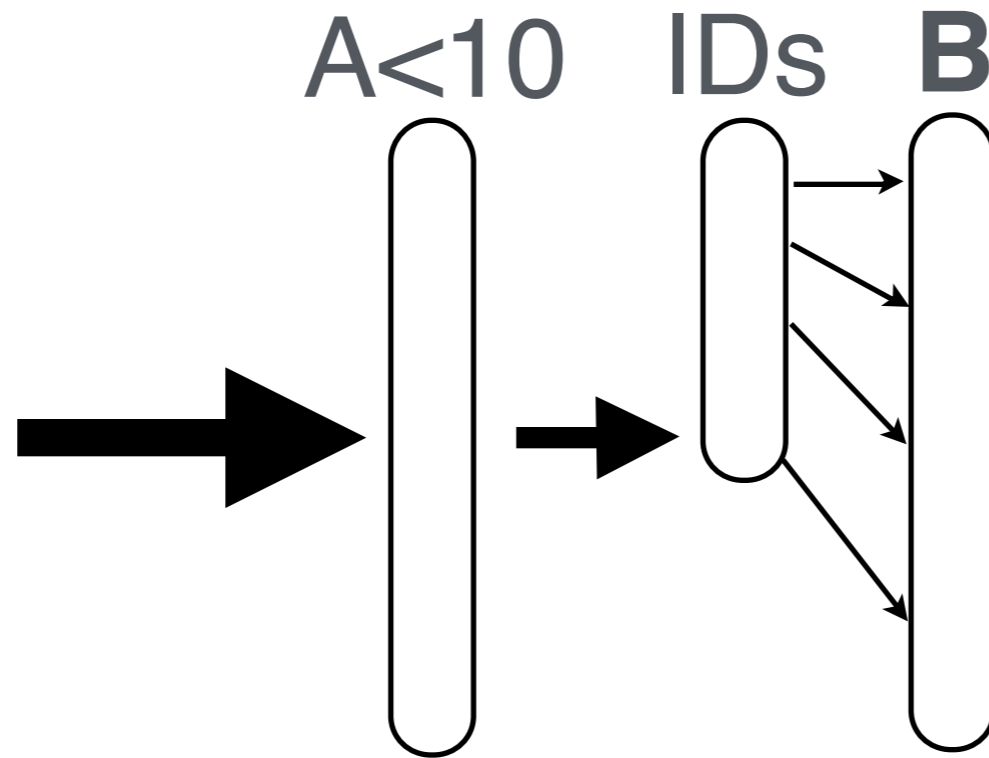
select min(C) from R where $A < 10$ & $B < 20$



disk



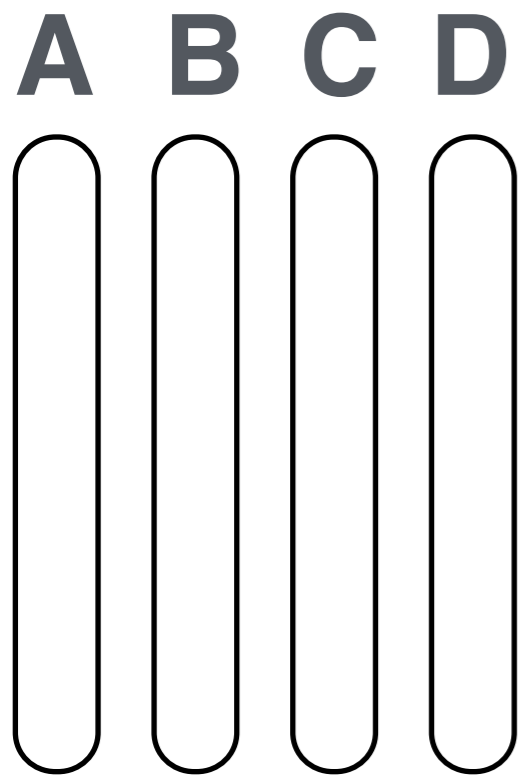
memory



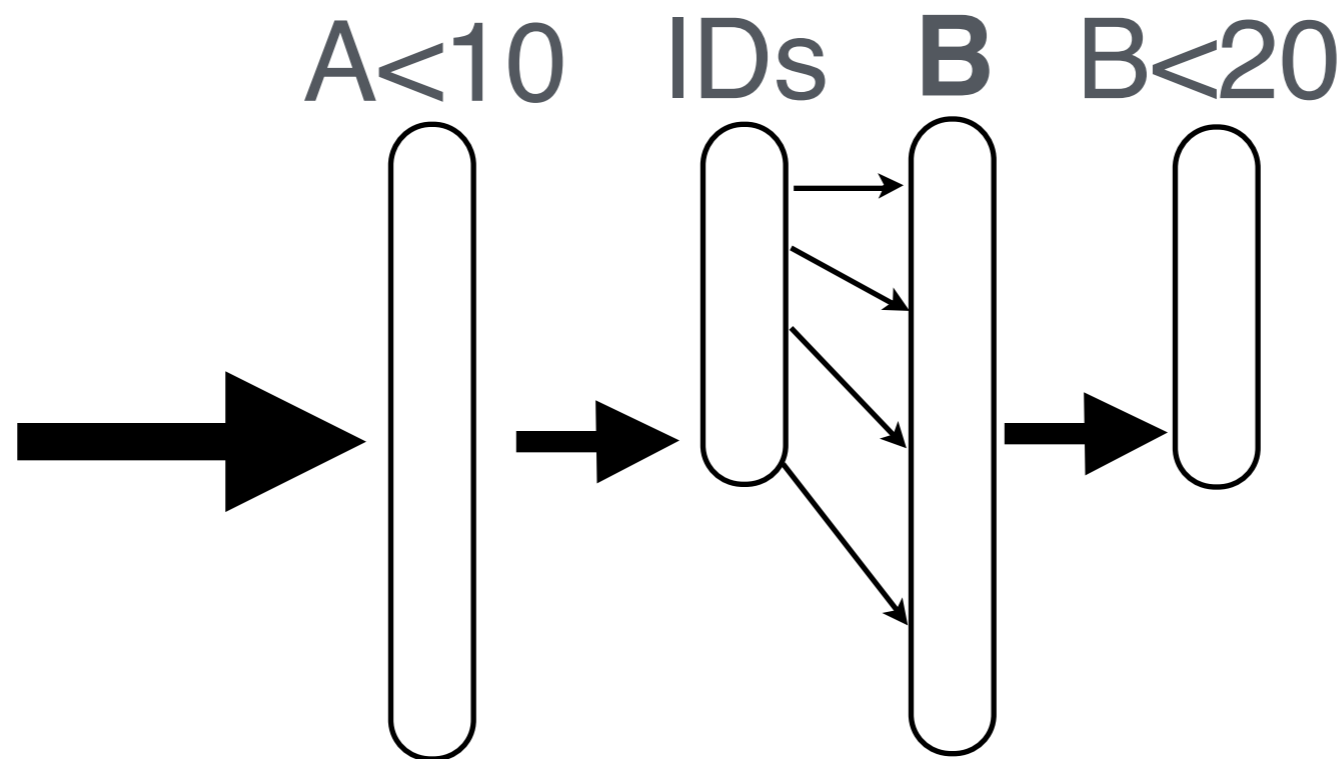
select min(C) from R where $A < 10$ & $B < 20$



disk



memory

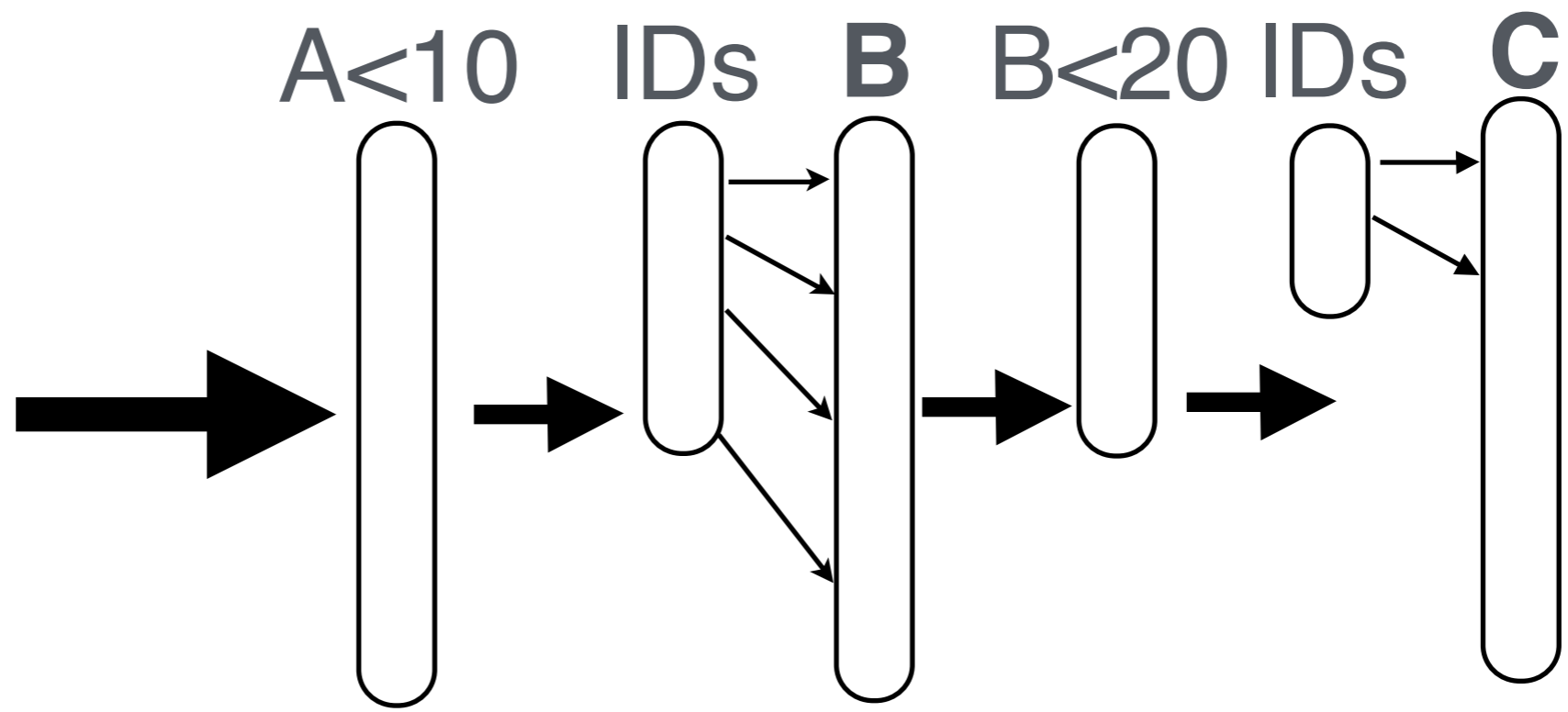
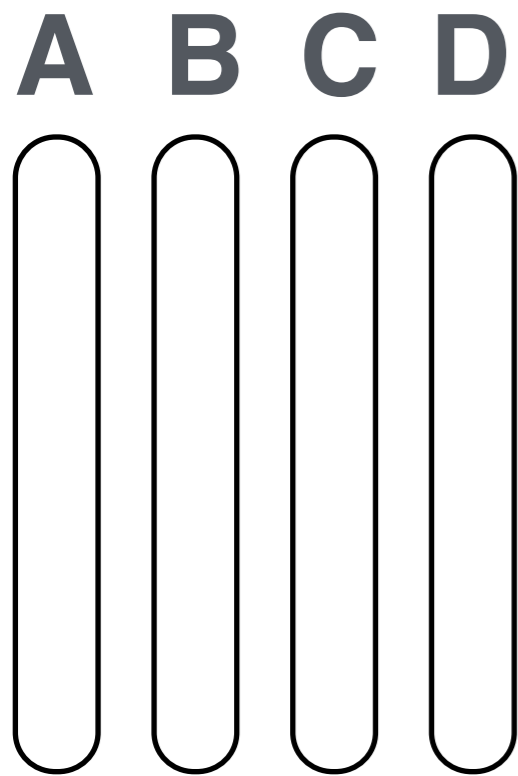


select min(C) from R where $A < 10$ & $B < 20$



disk

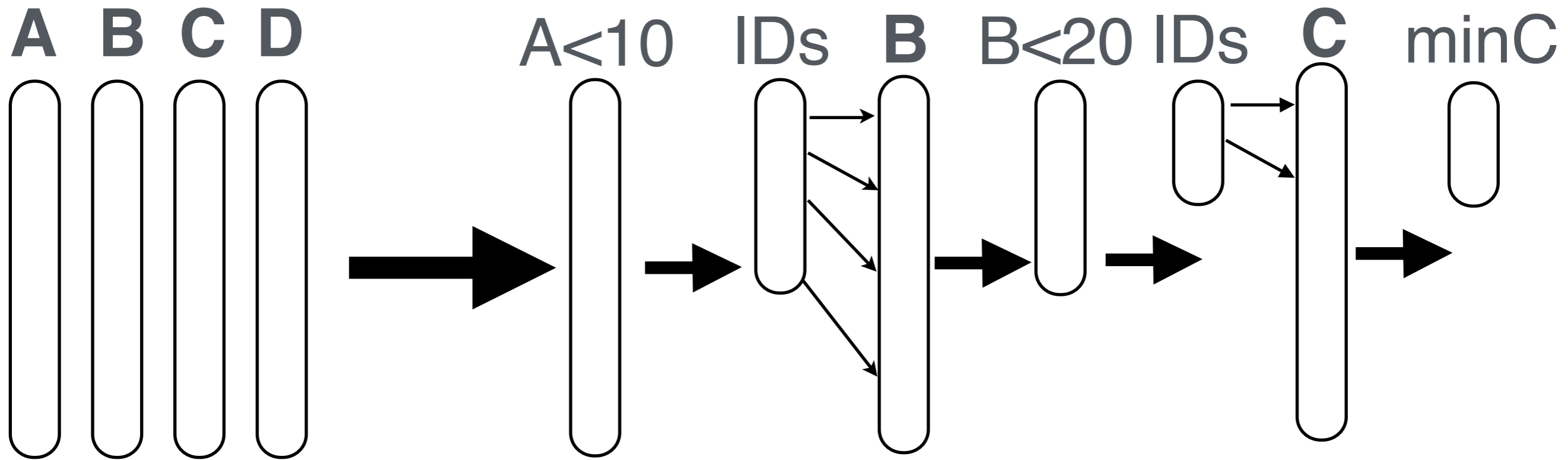
memory



select min(C) from R where $A < 10$ & $B < 20$

↓
disk

memory

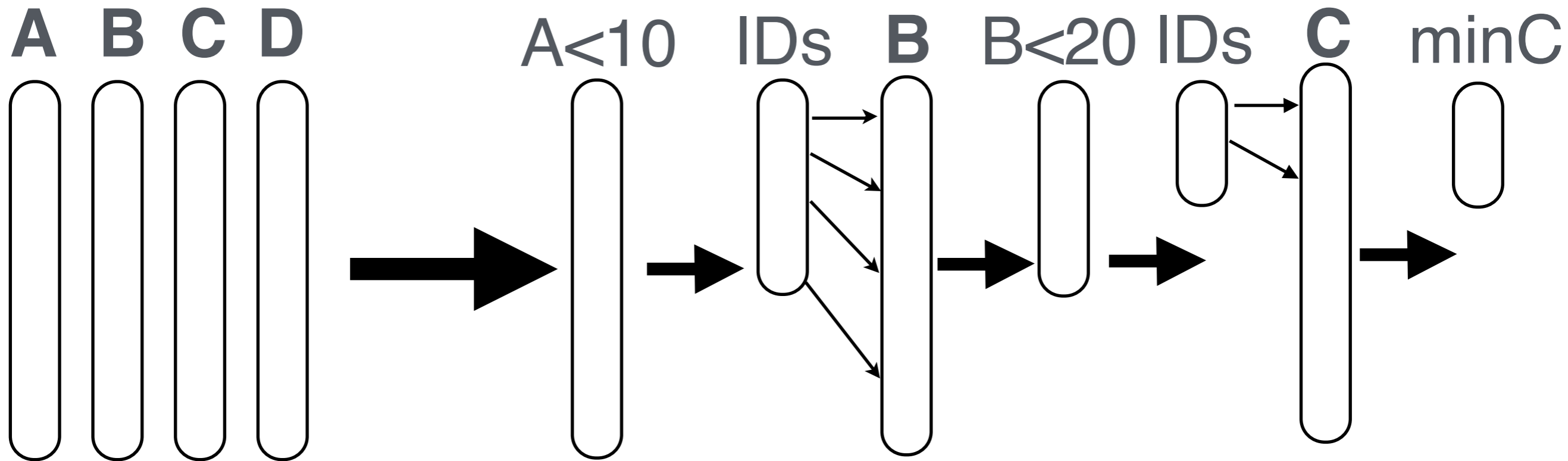


select min(C) from R where A < 10 & B < 20

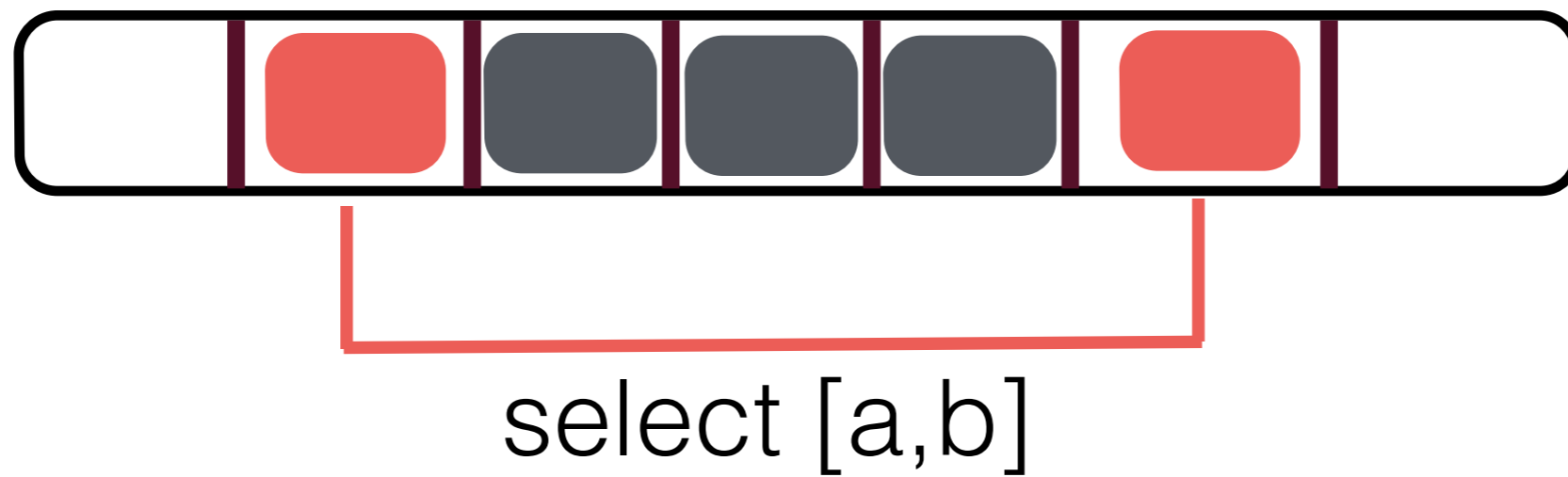


disk

memory



column lock and release
as soon as an operator completes

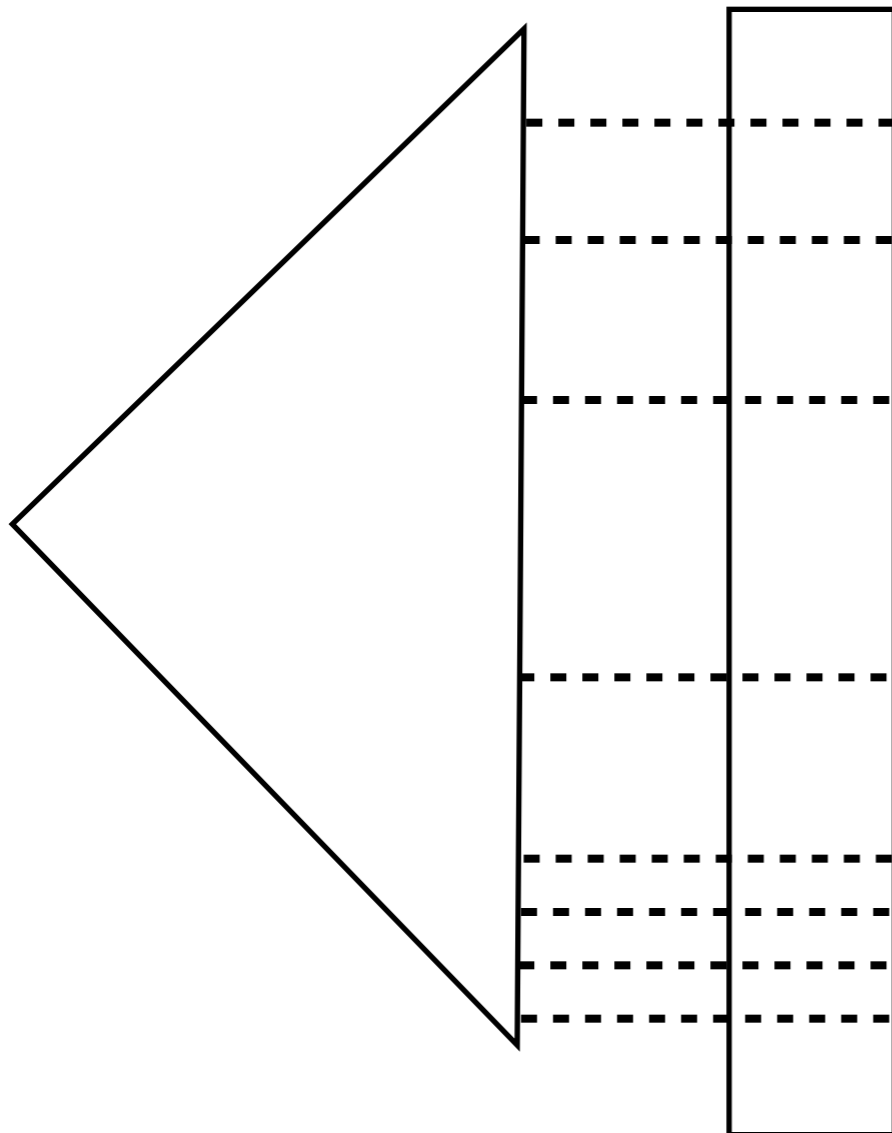


need to latch only to be cracked
pieces (max 2 per select)

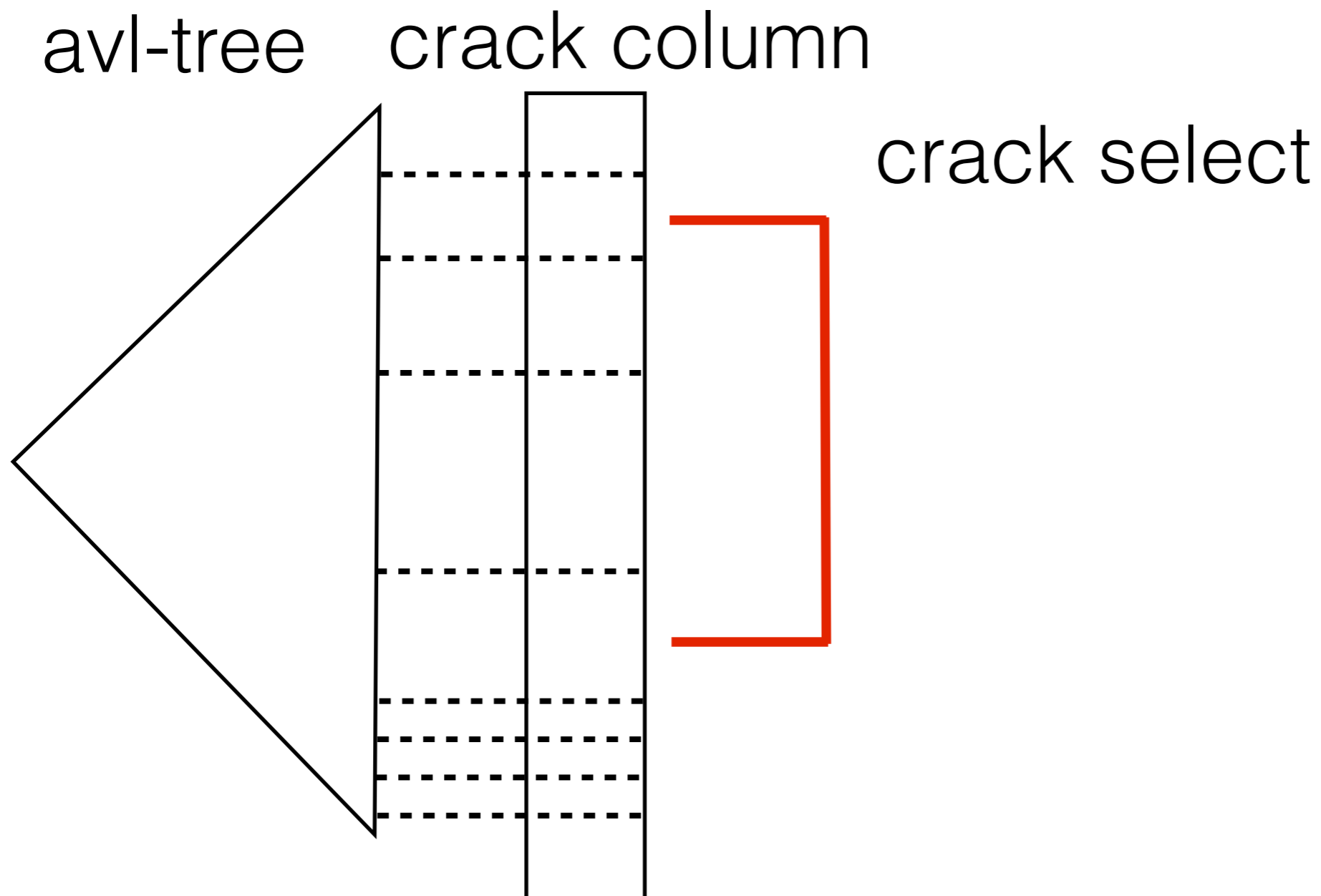


piece locking

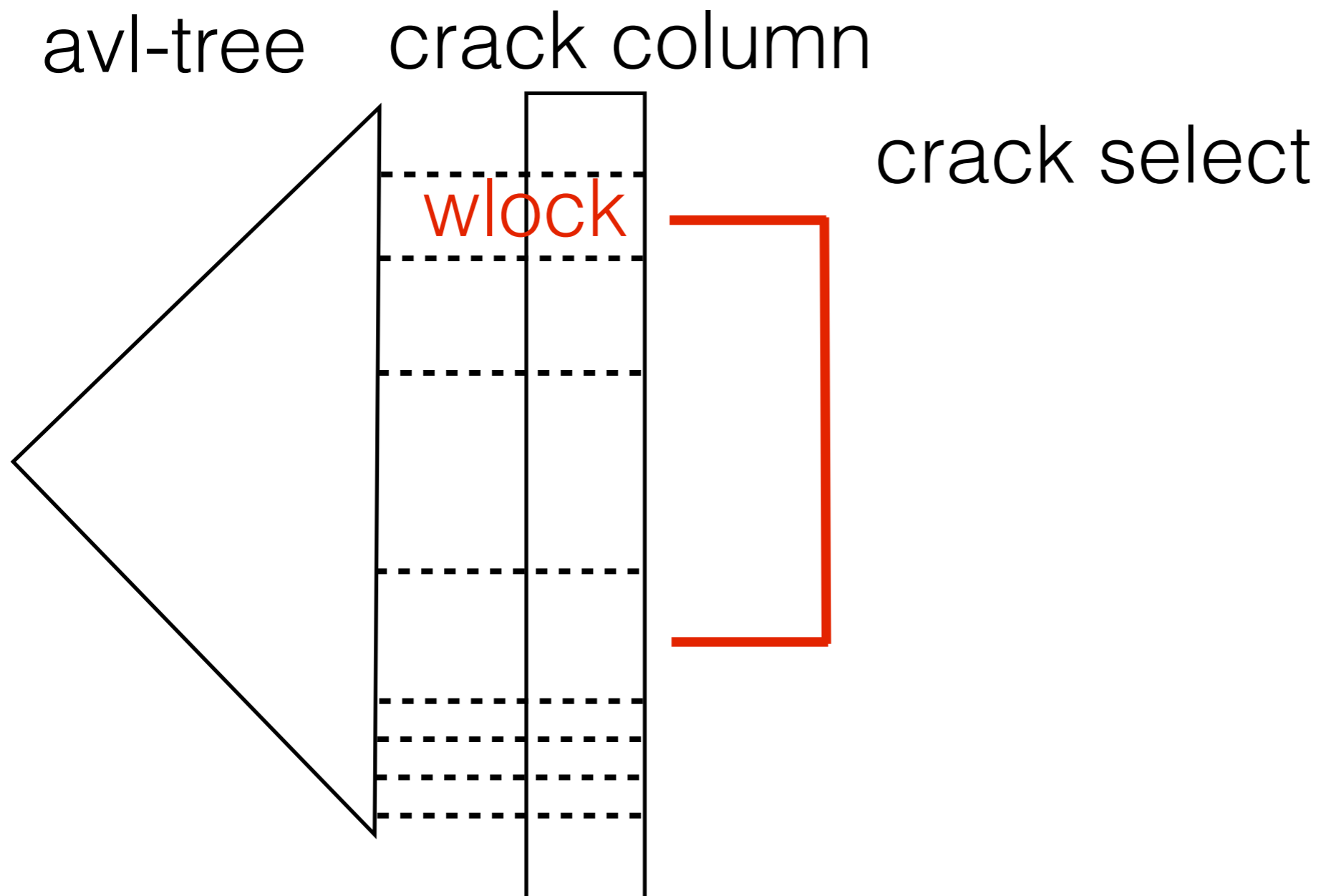
avl-tree crack column



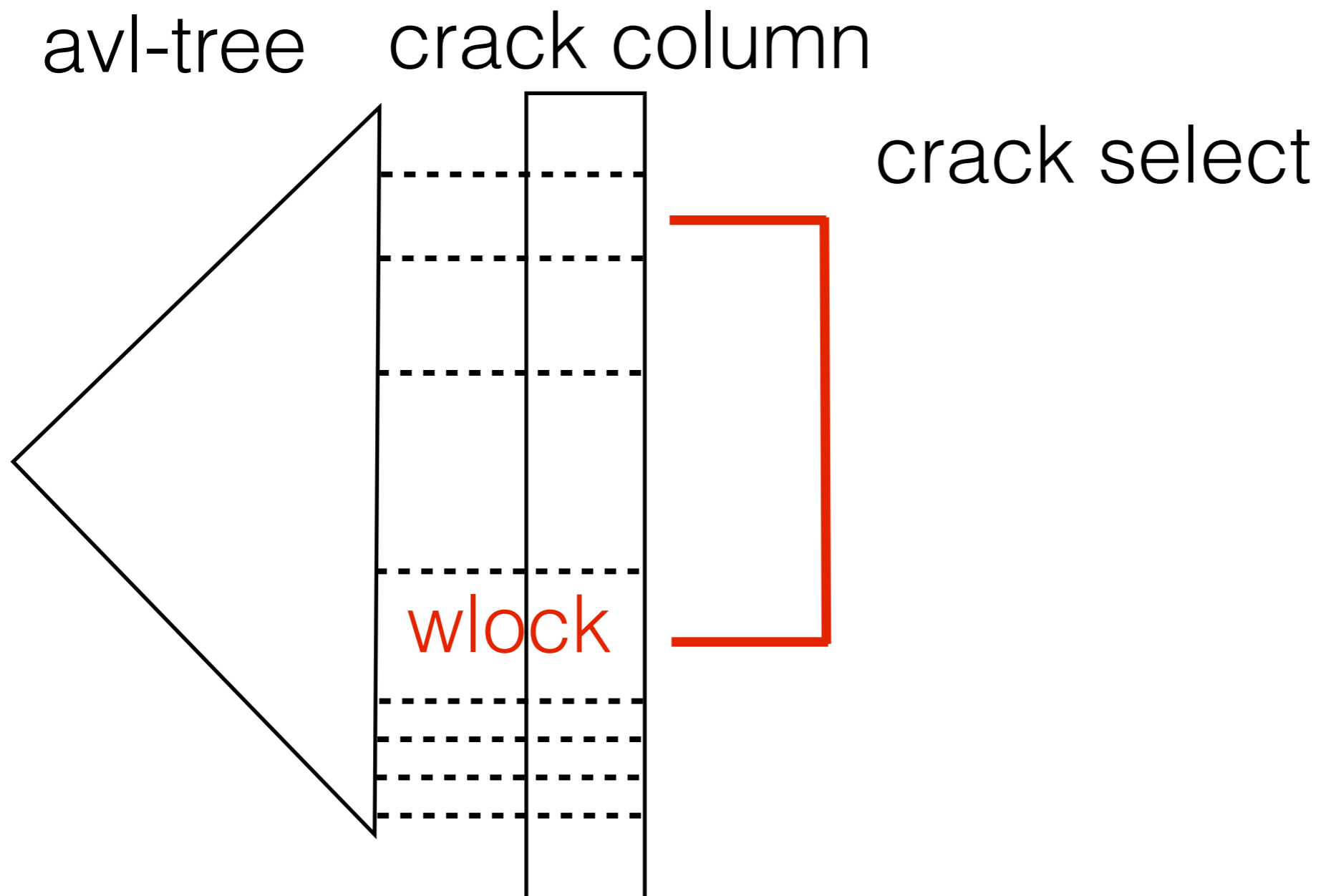
piece locking



piece locking

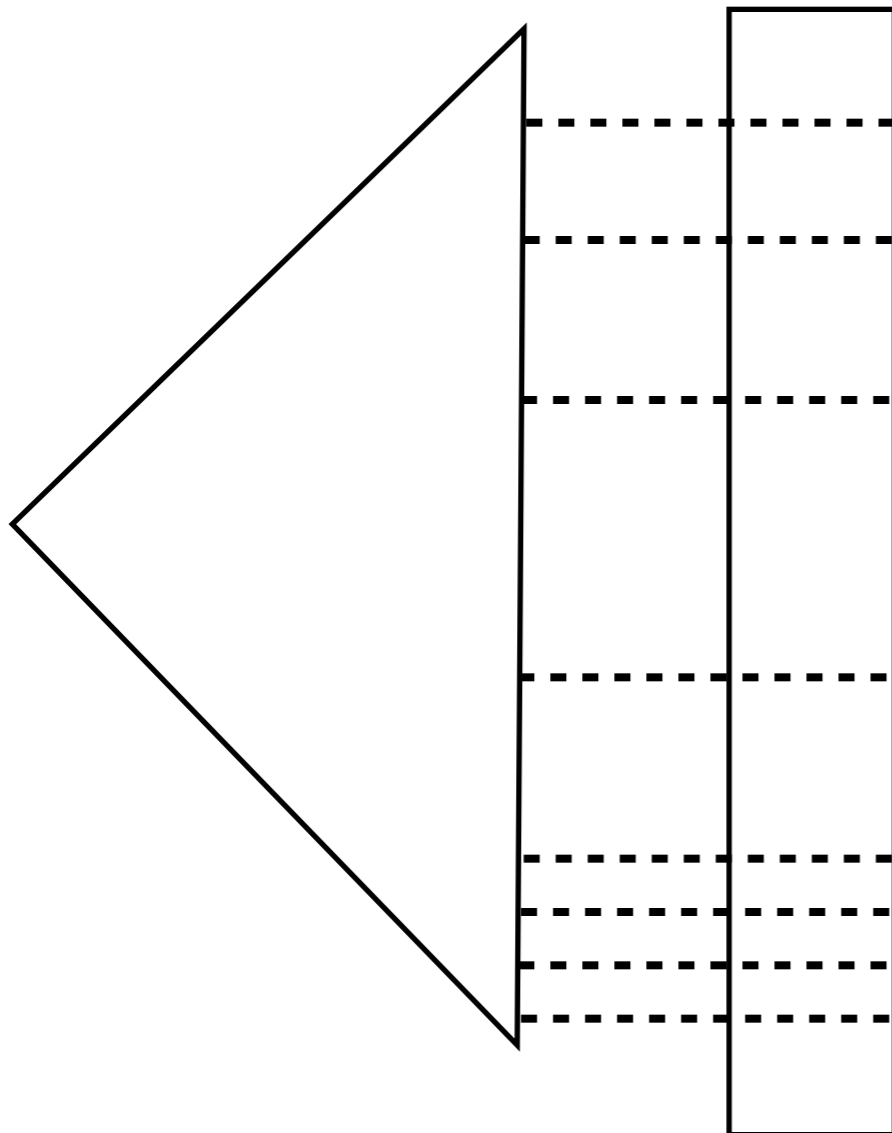


piece locking

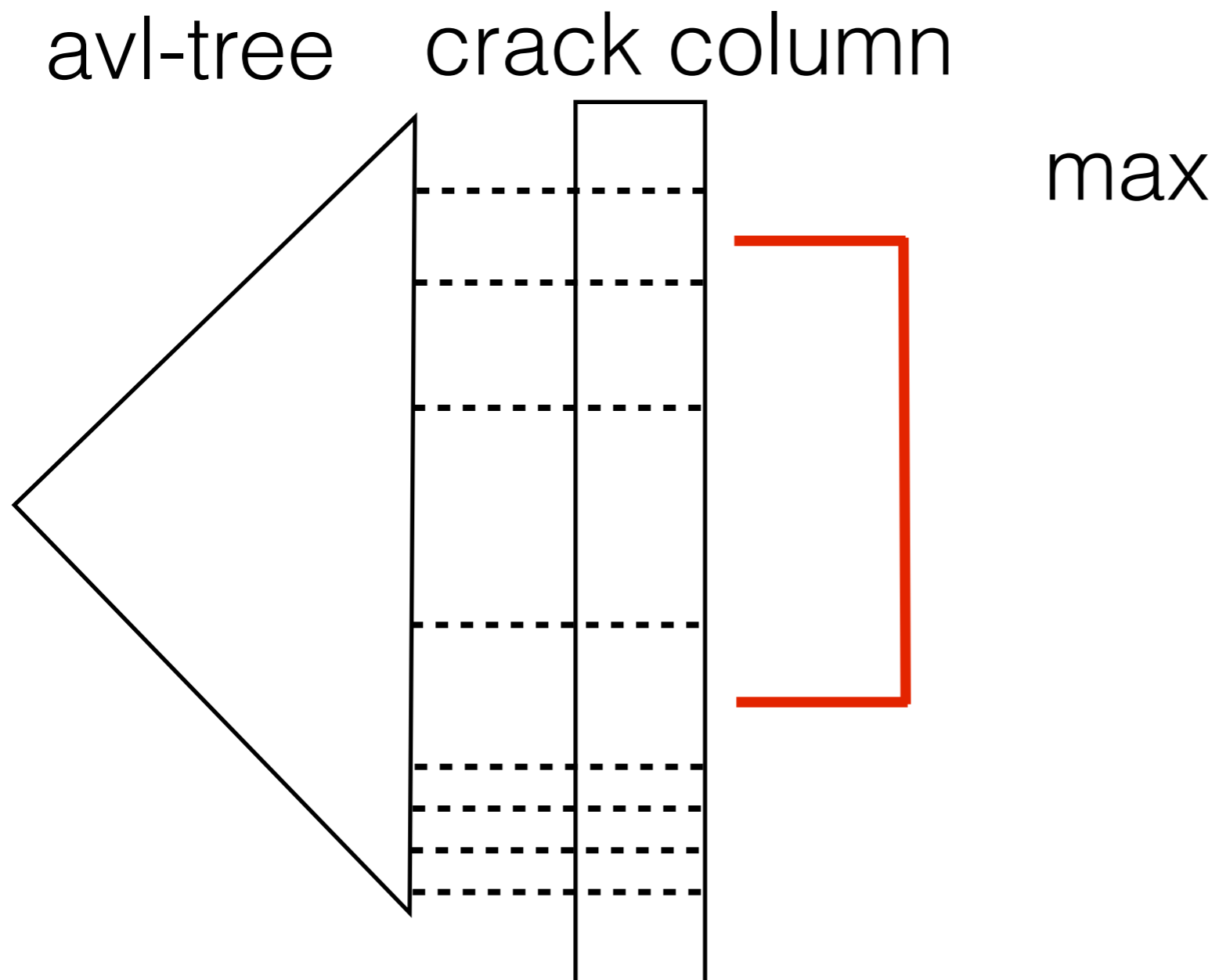


piece locking

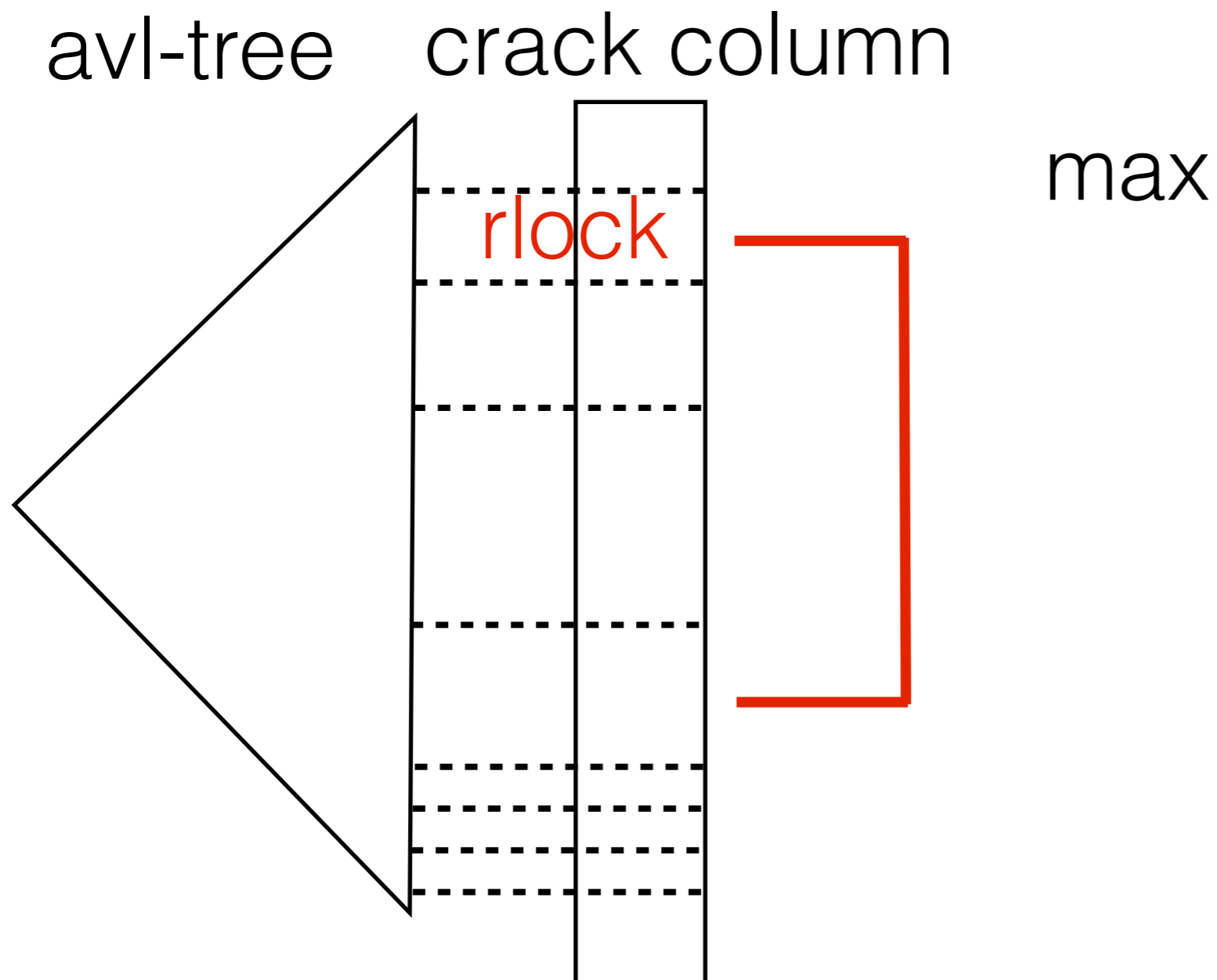
avl-tree crack column



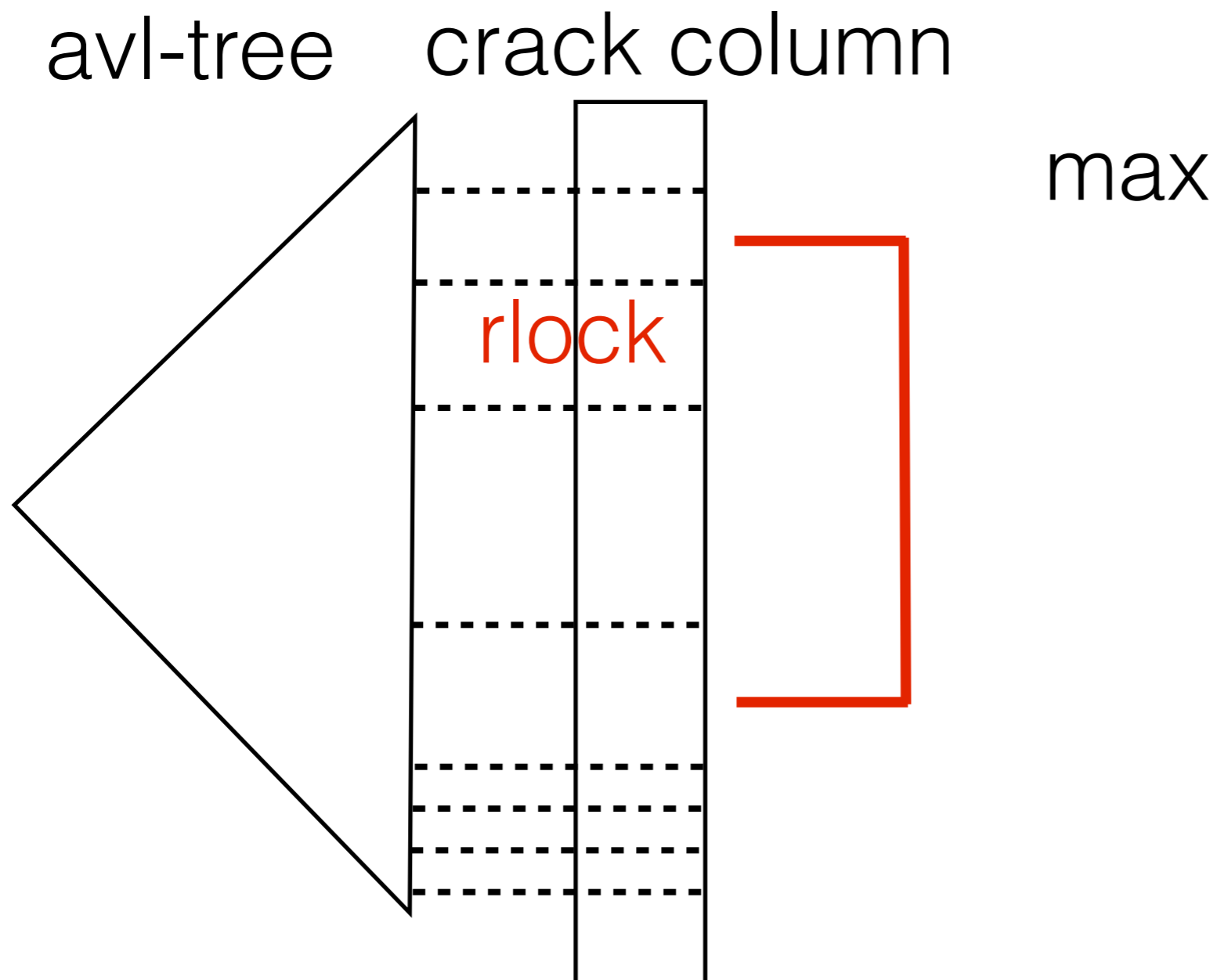
piece locking



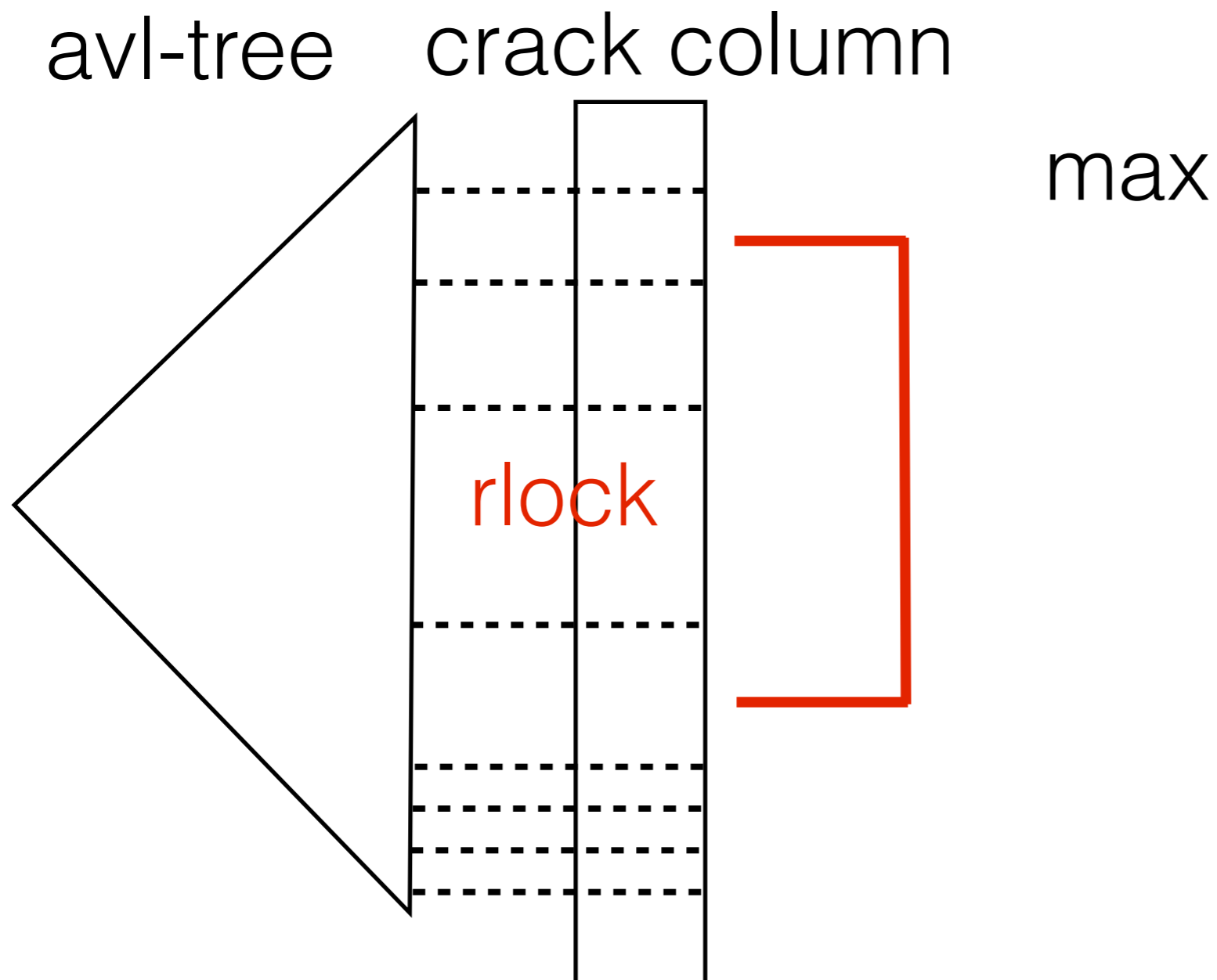
piece locking



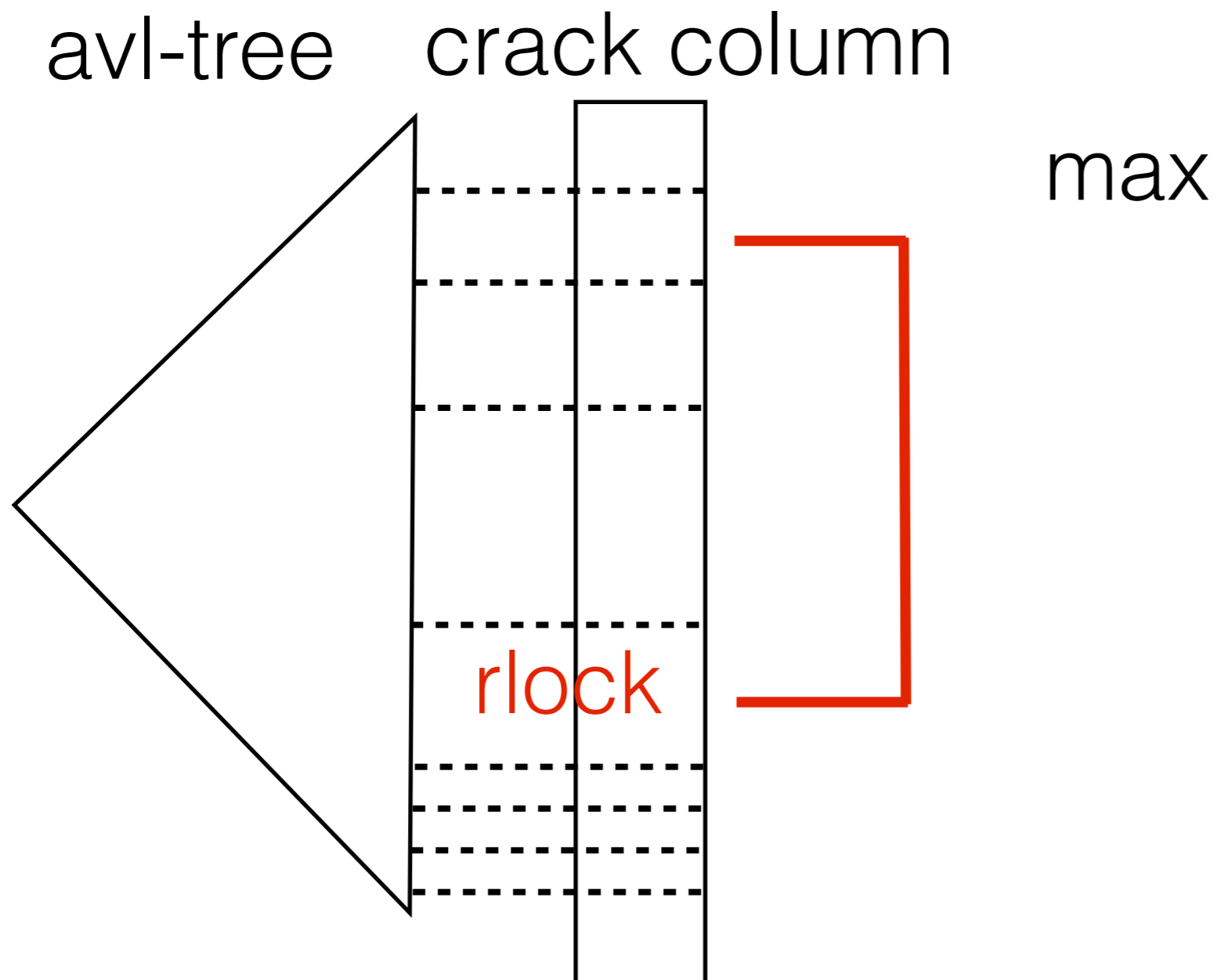
piece locking



piece locking

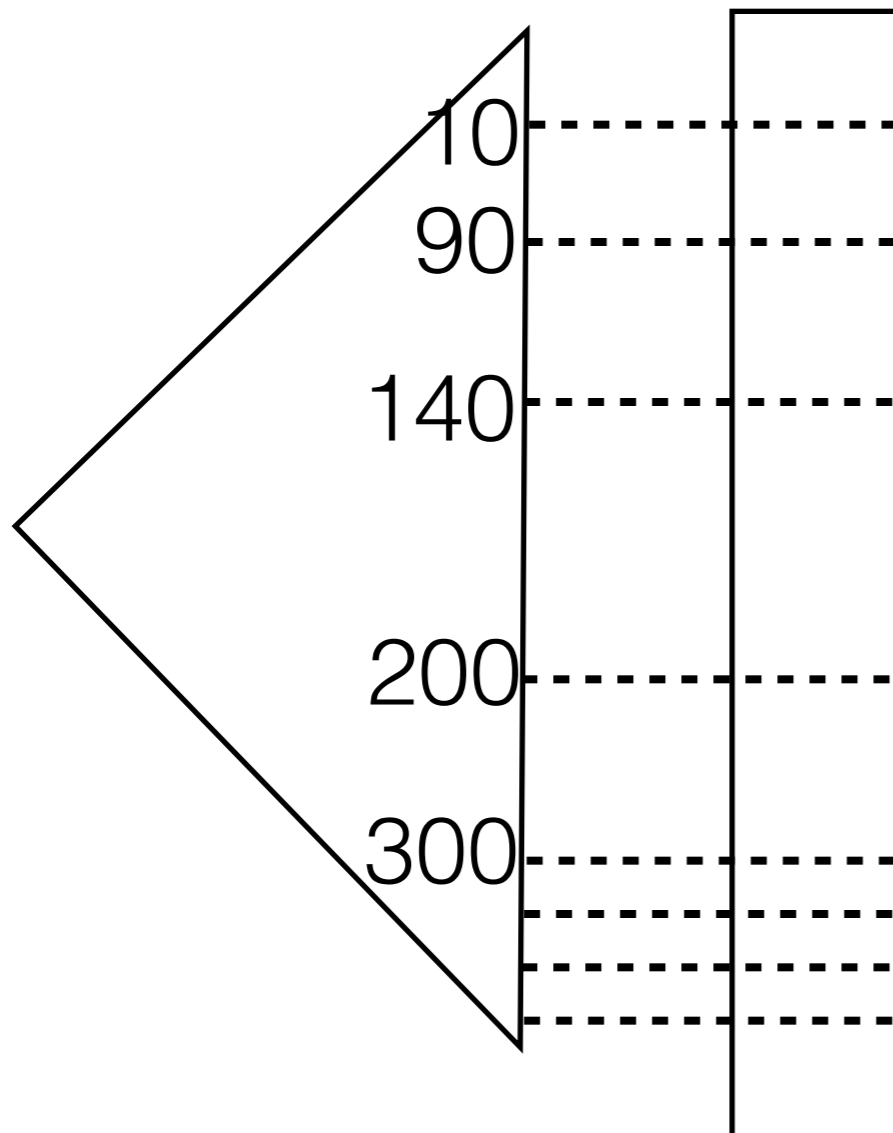


piece locking



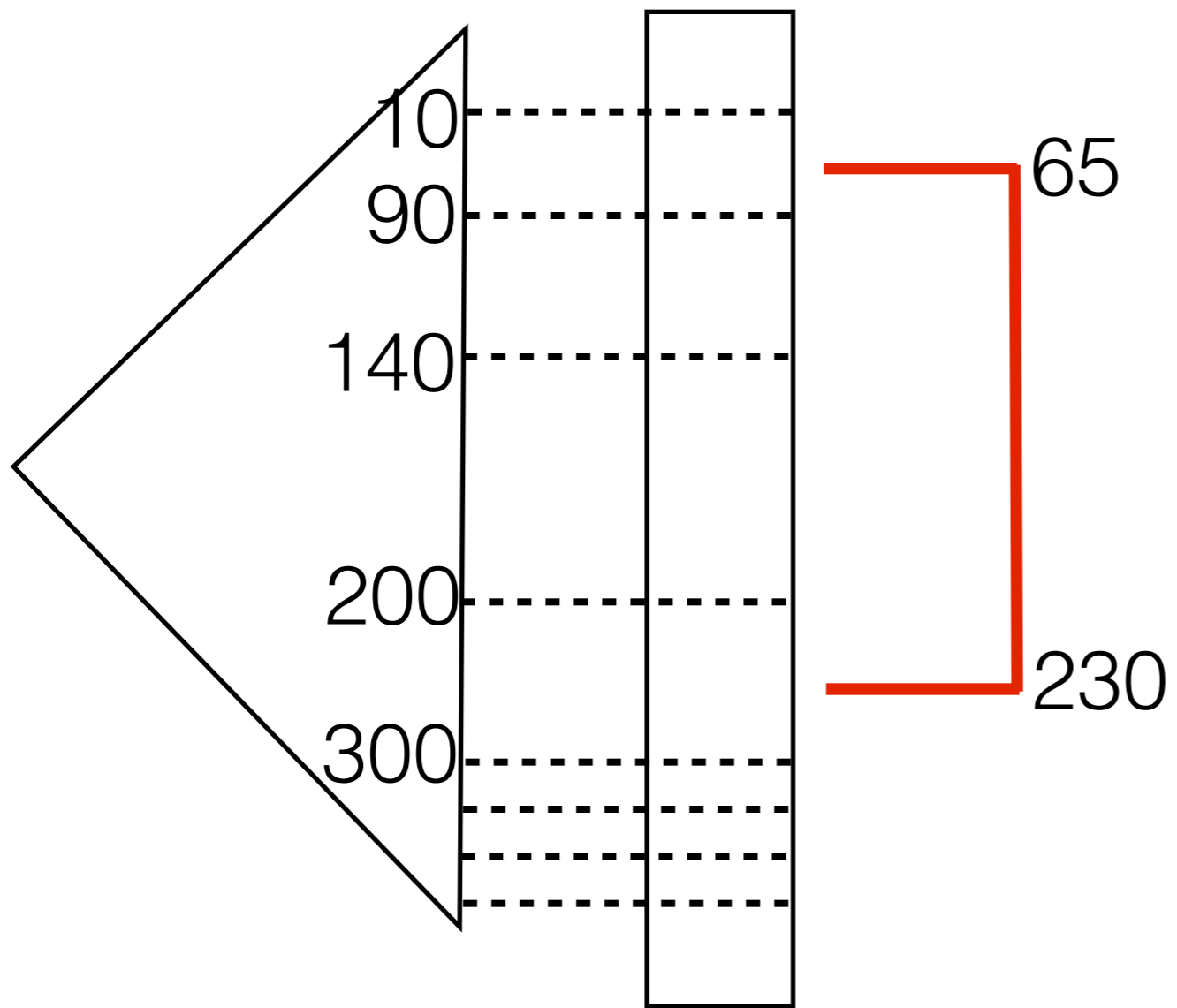
piece locking

avl-tree crack column



piece locking

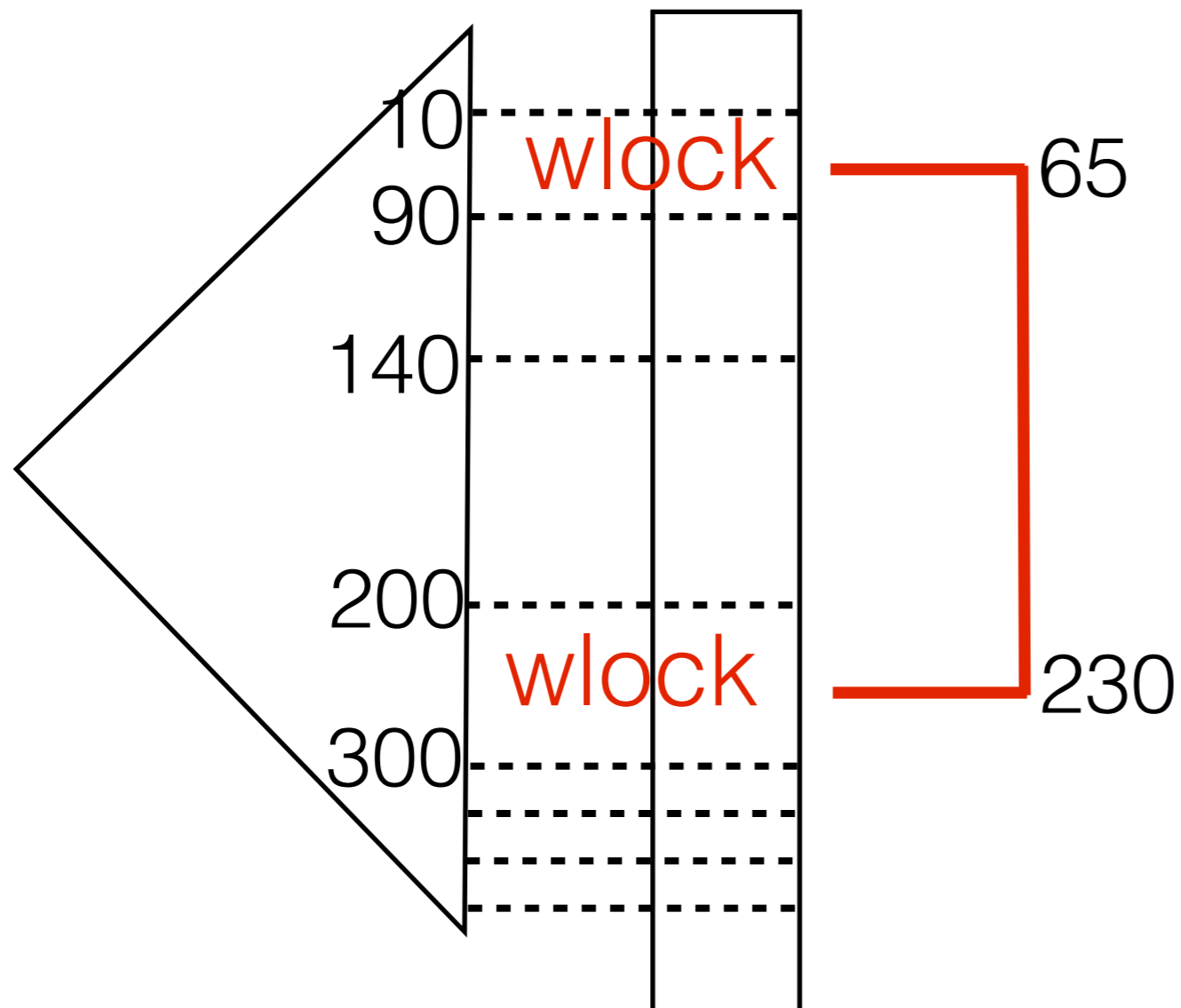
avl-tree crack column



crack select

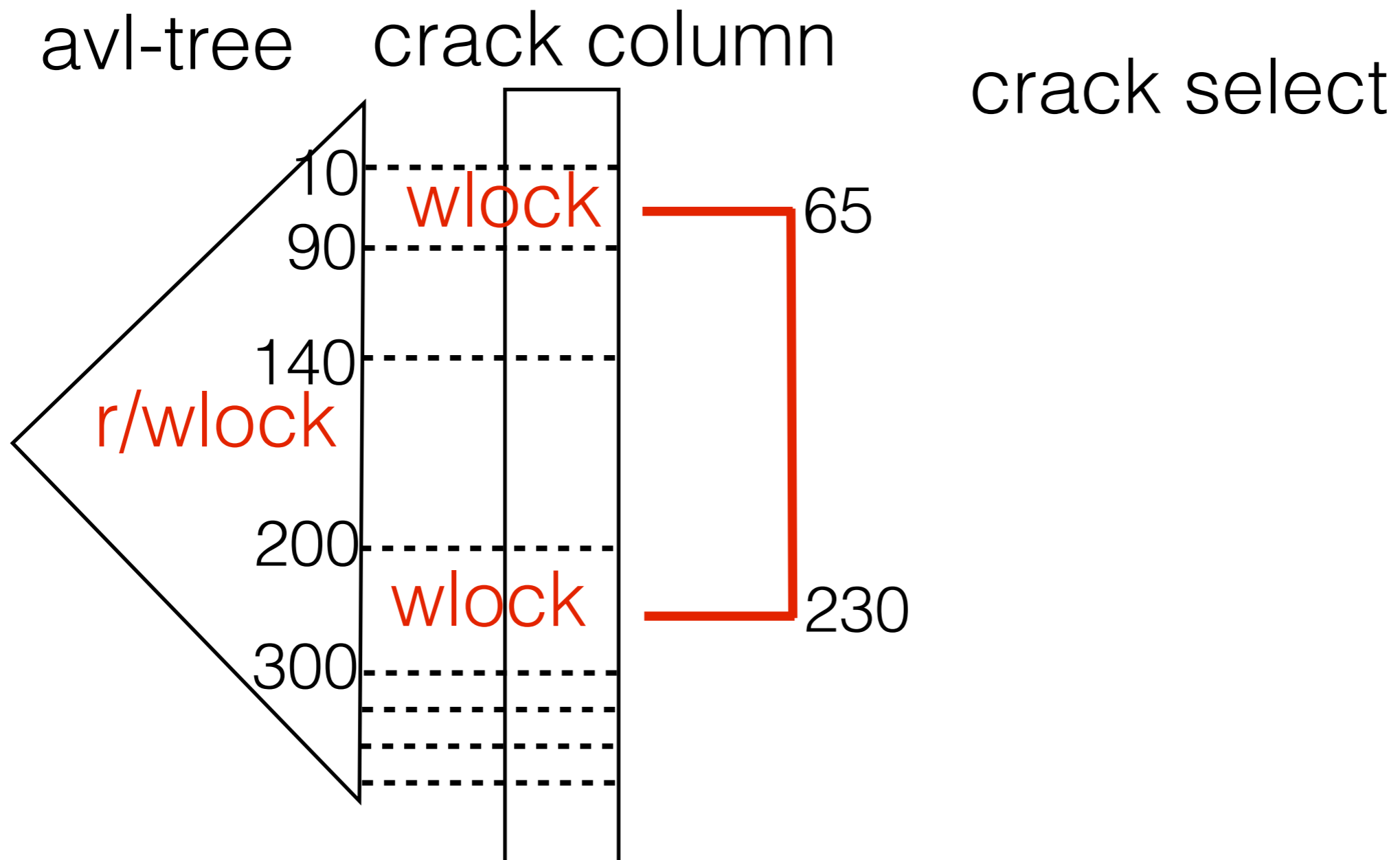
piece locking

avl-tree crack column

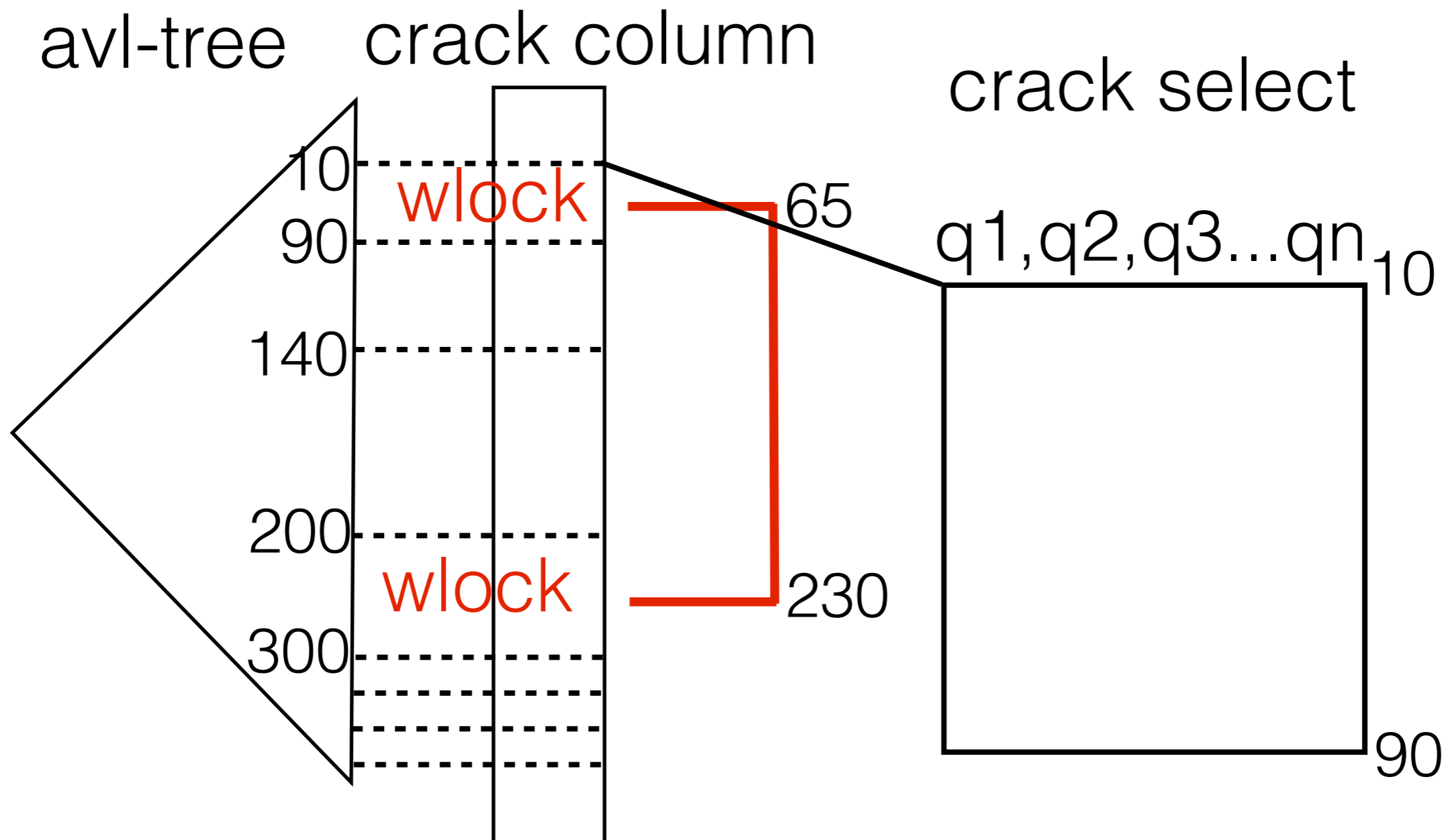


crack select

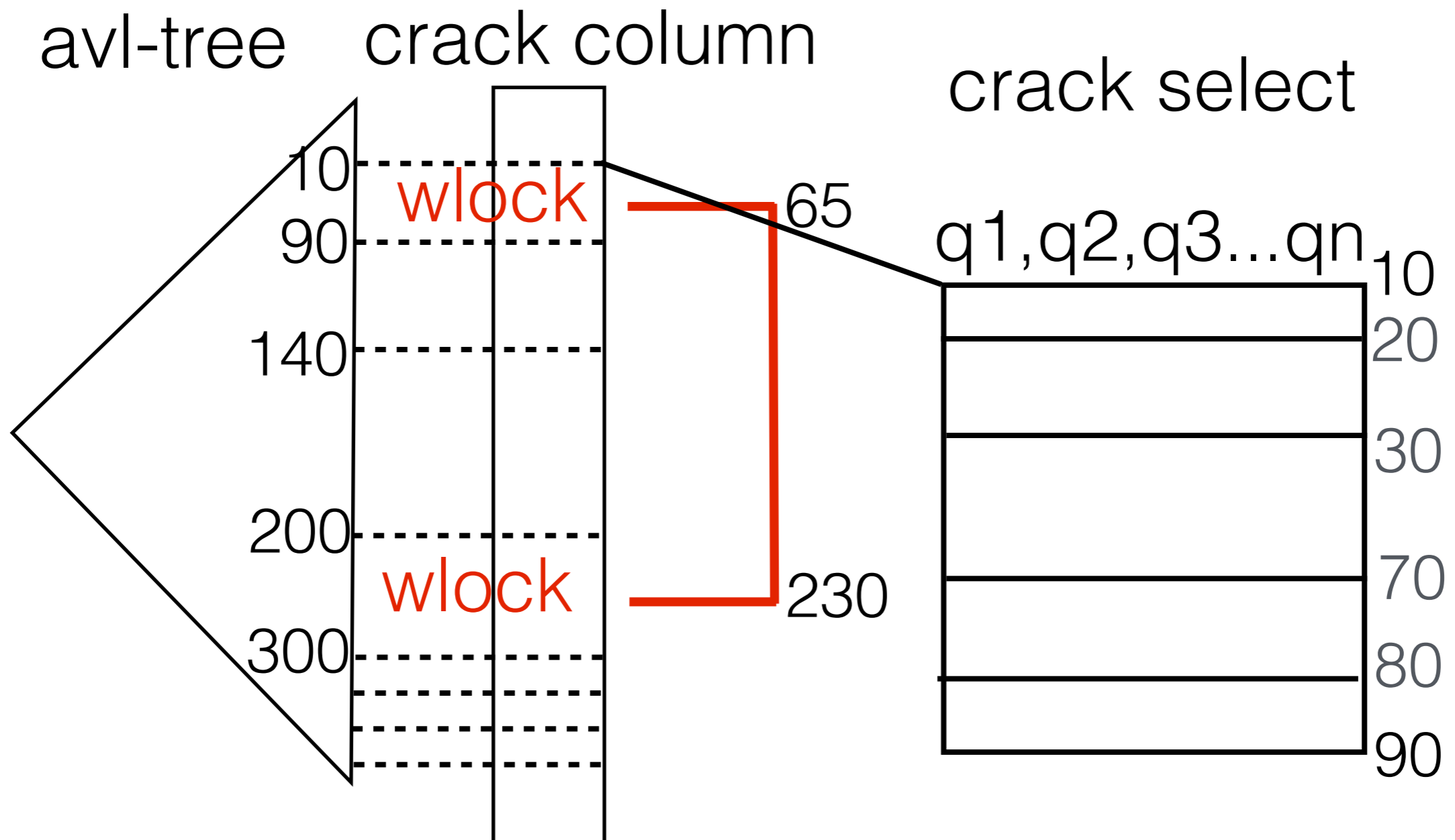
piece locking



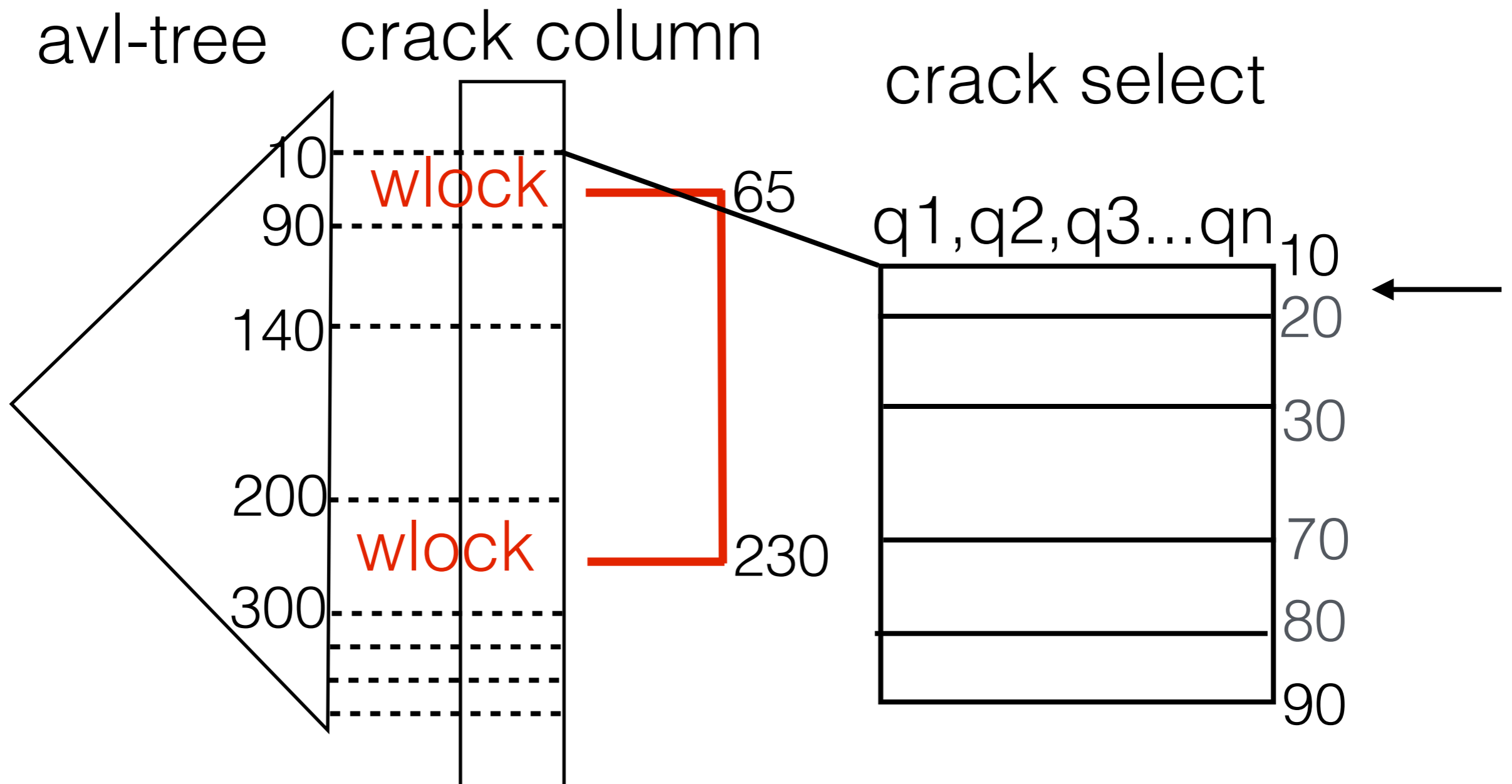
piece locking



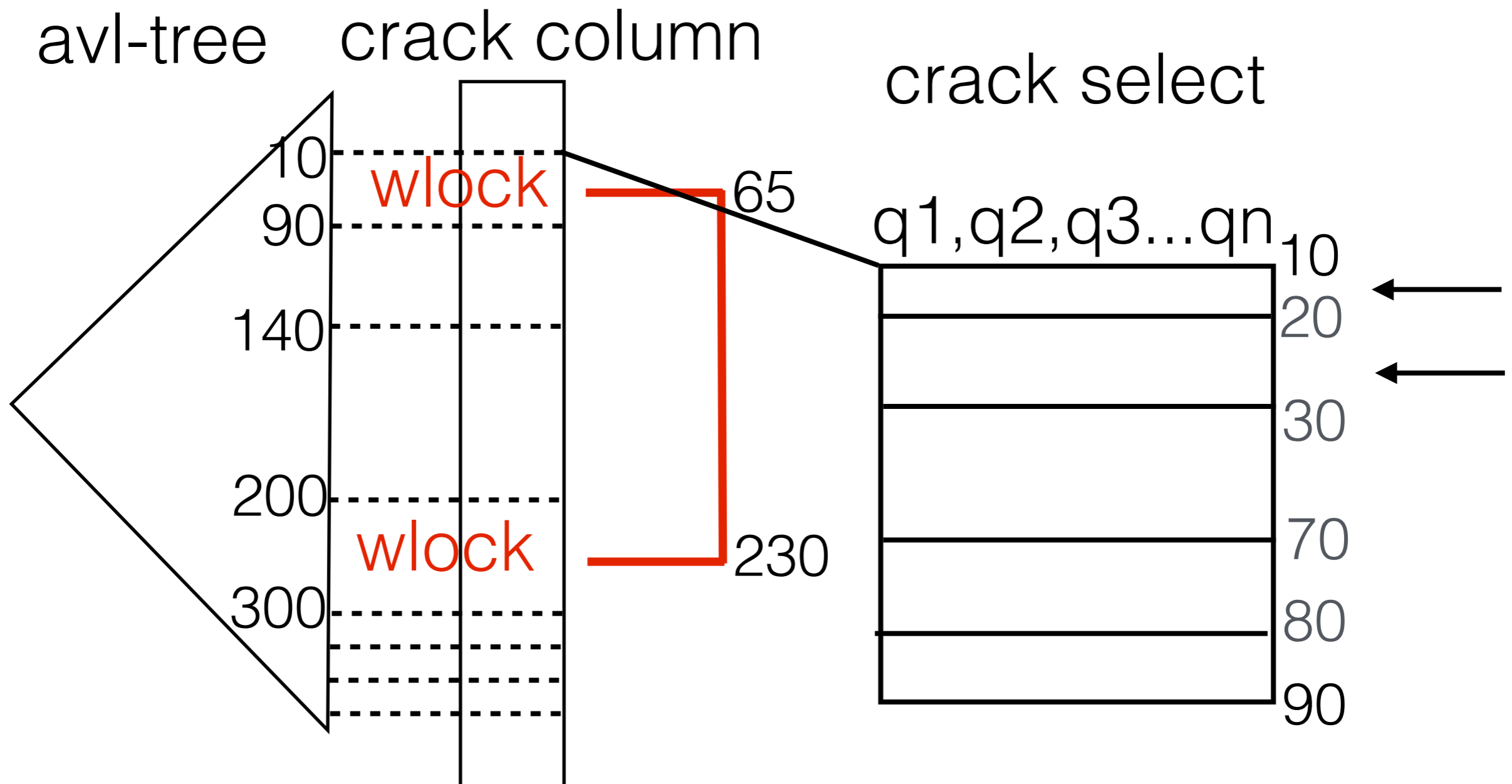
piece locking



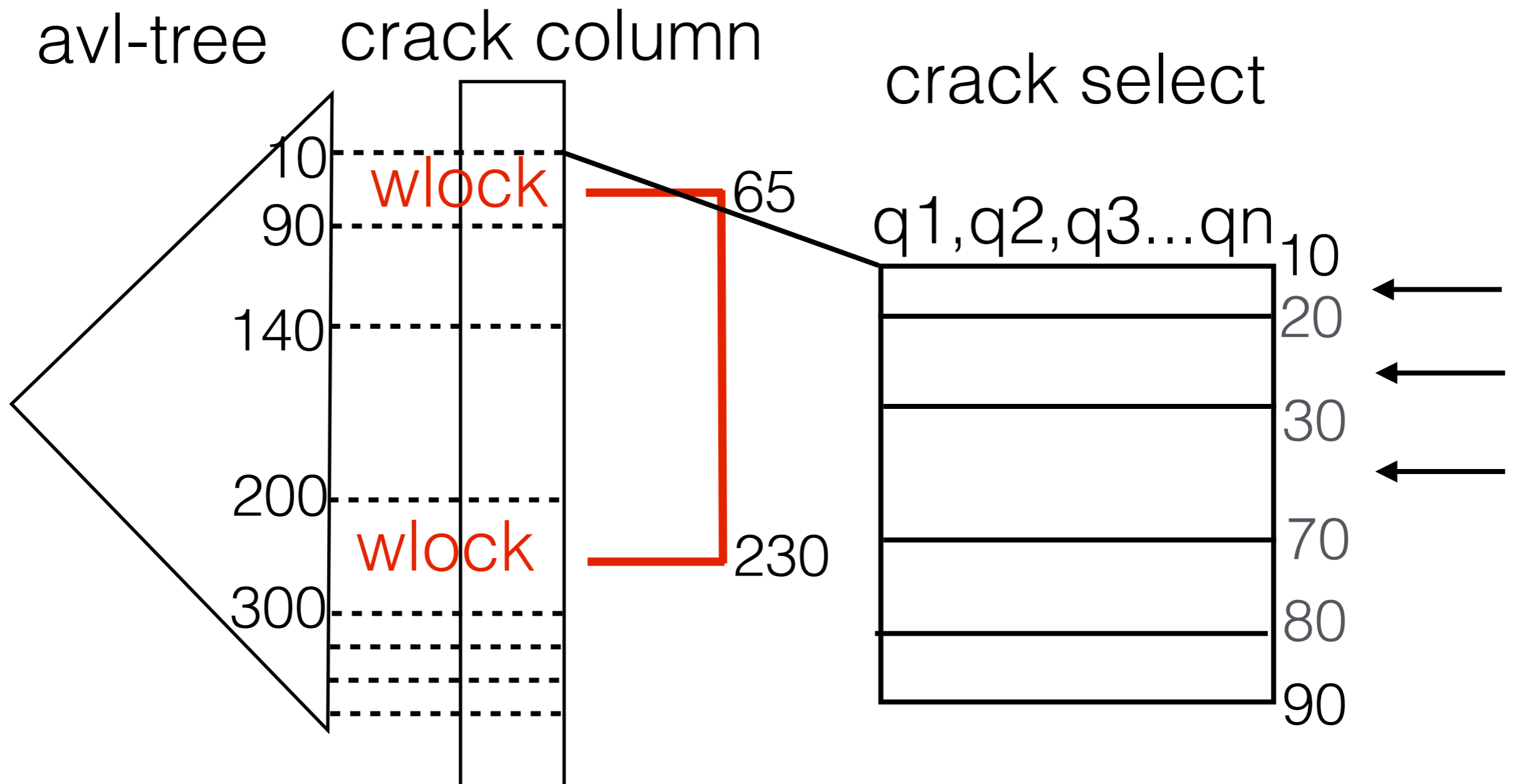
piece locking



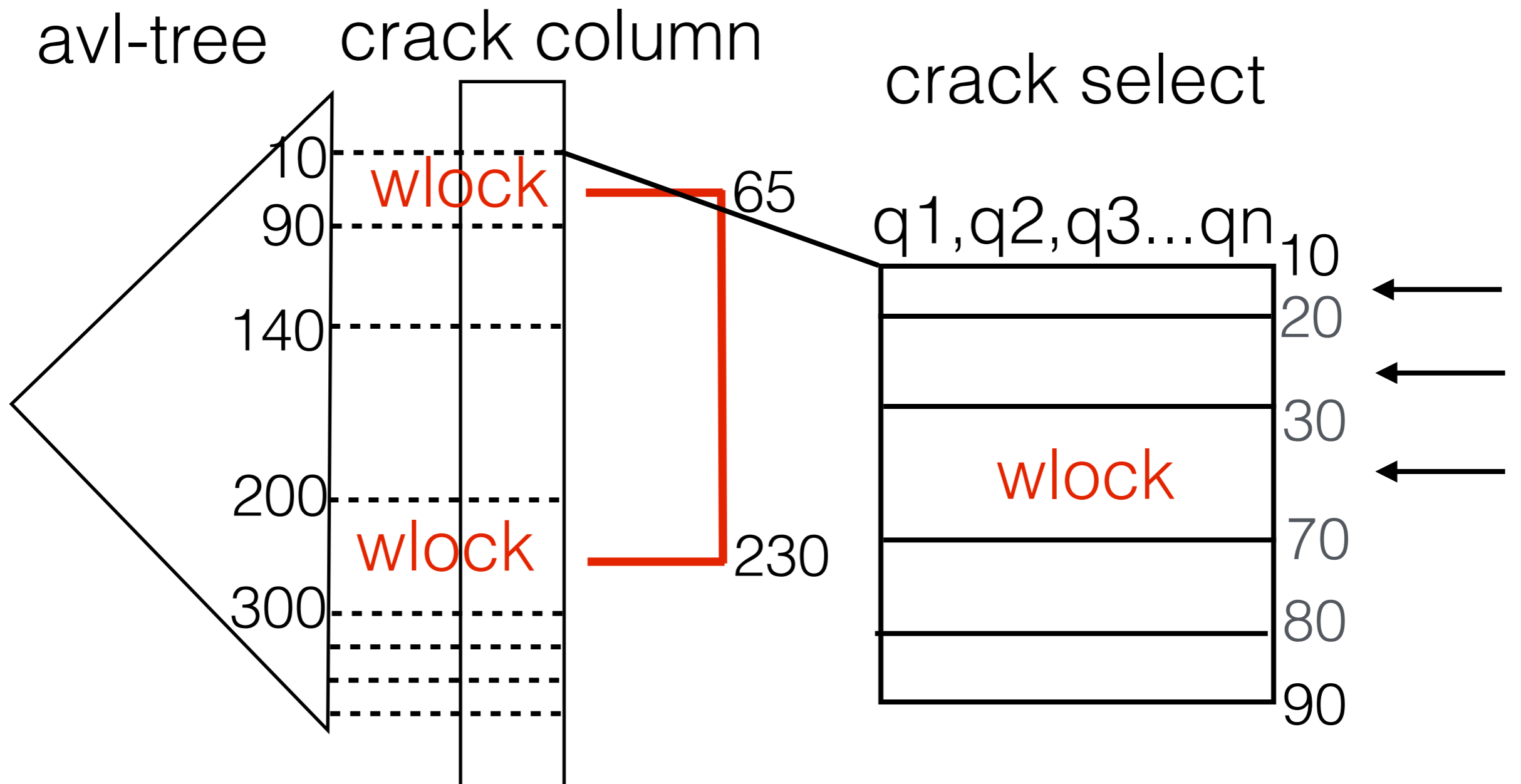
piece locking



piece locking



piece locking



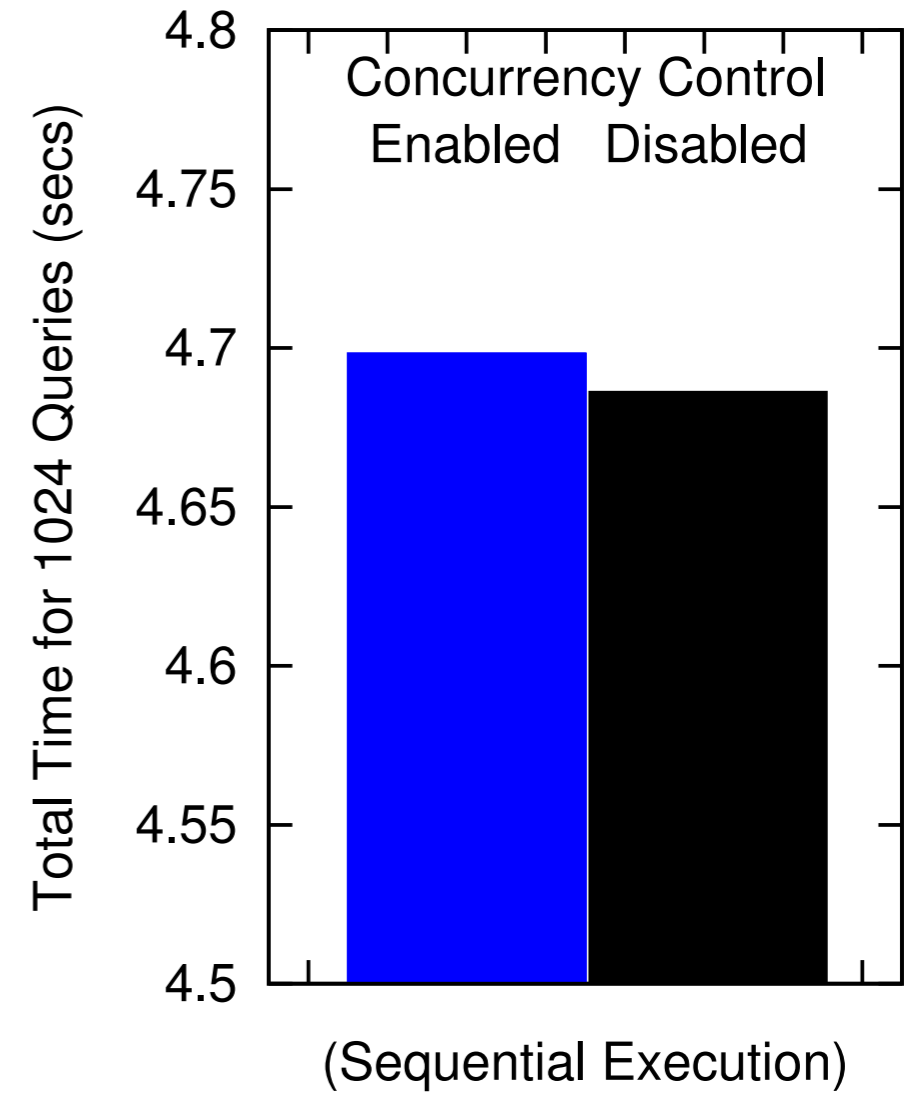
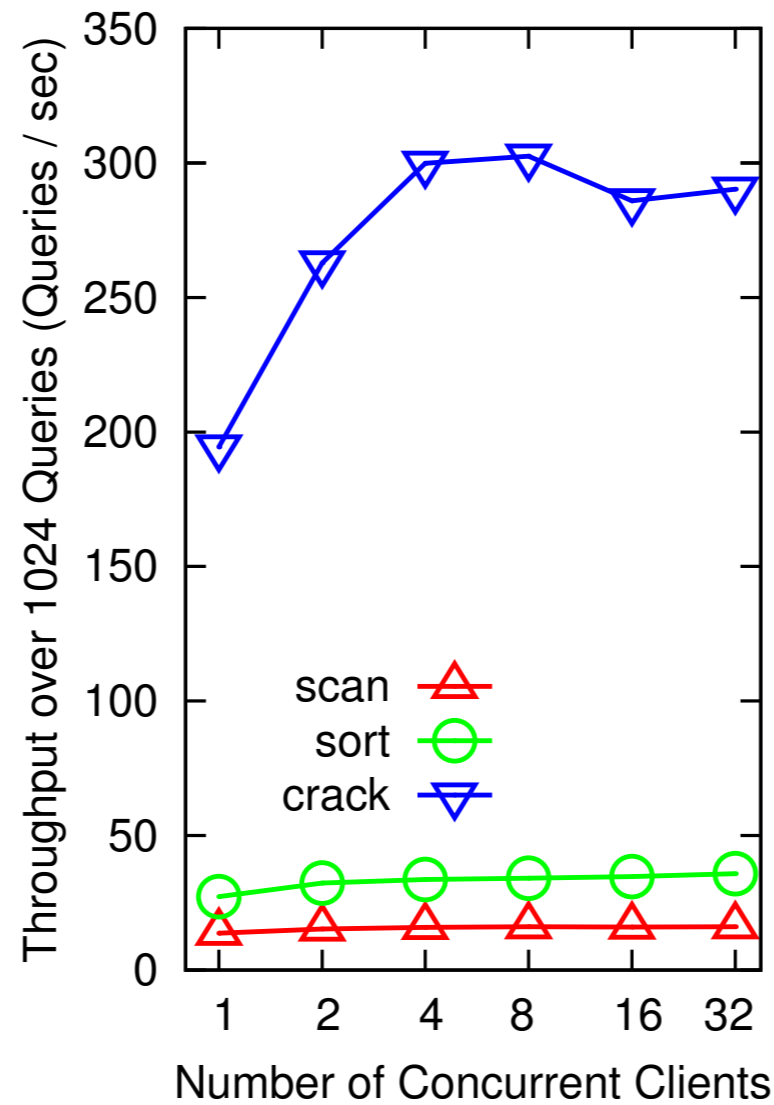
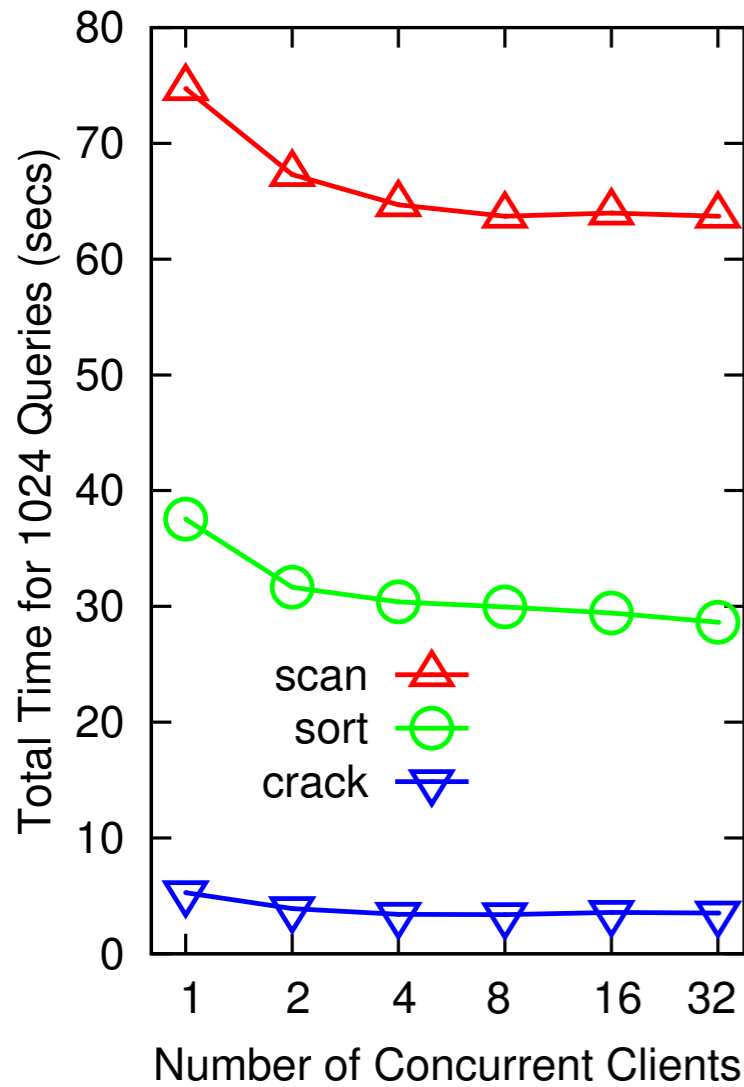
Rows: 100M

Query: select sum(A) from R where v1 < A1 < v2

Selectivity: 0.01%, Random, # of queries: 1024

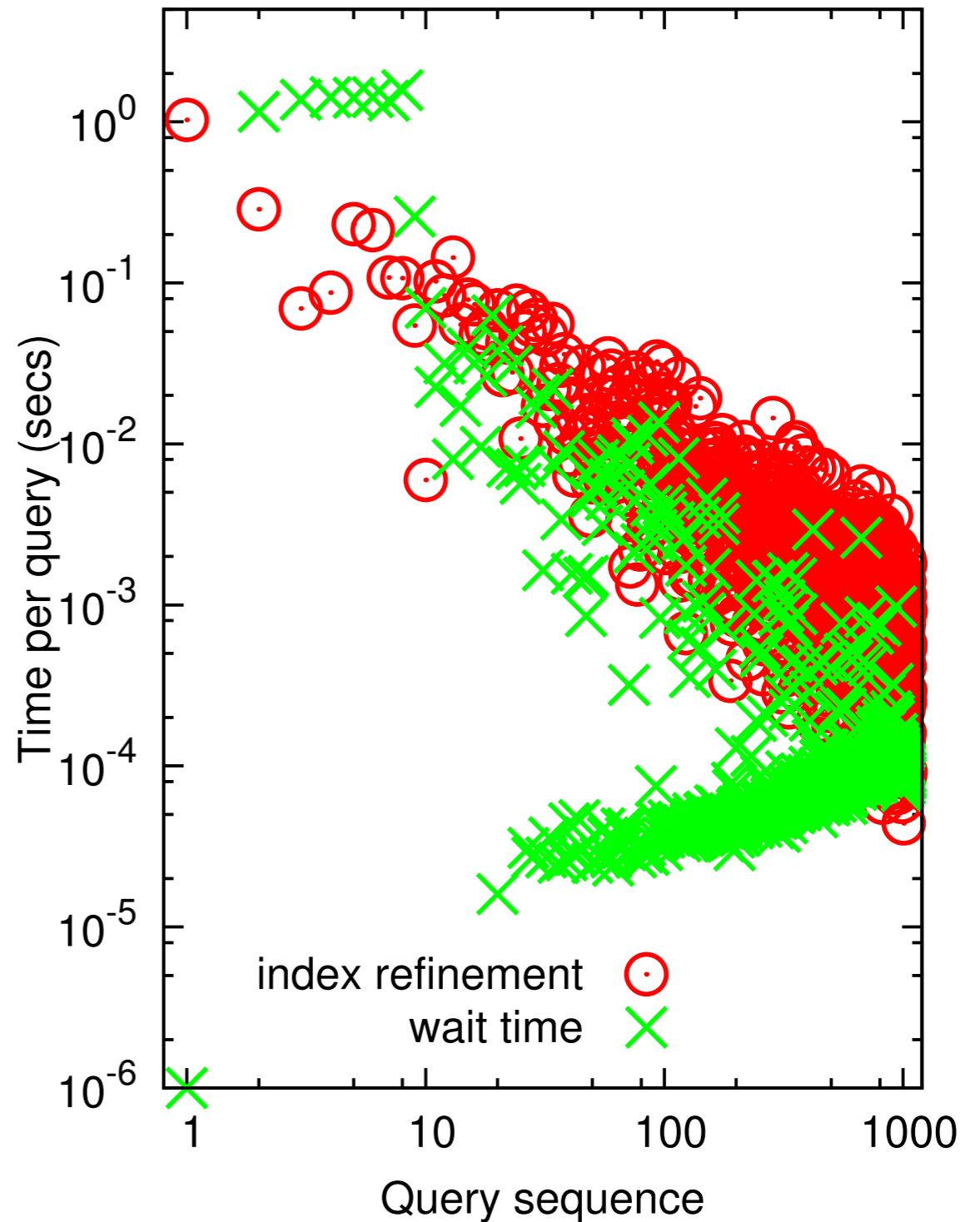
Clients: 1-32, Machine: 4 cores

adaptive indexing maintains its performance advantage



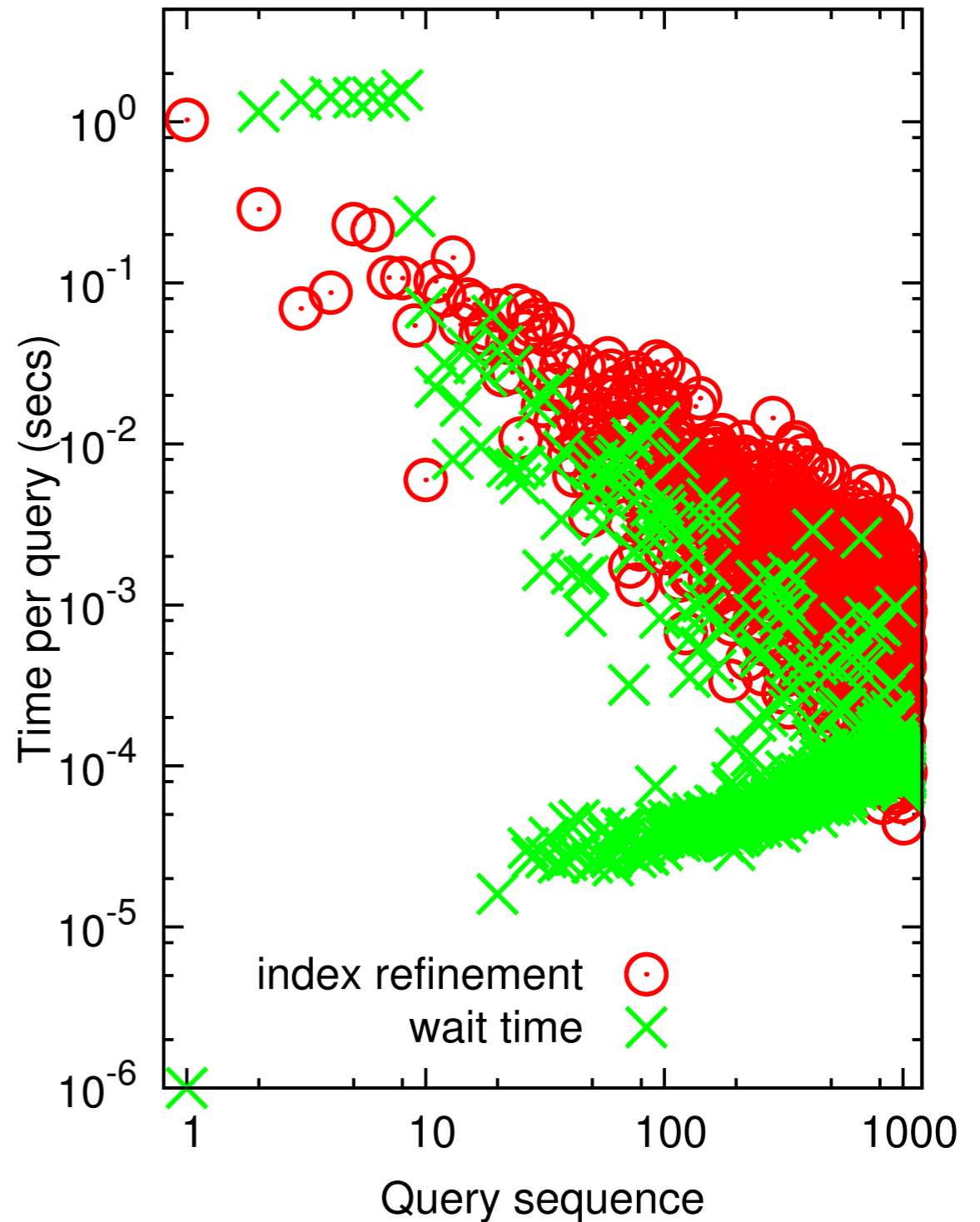
crack time and wait time

Query: select sum(A) from R where v1 < A1 < v2
Selectivity: 0.01%, Random, # of queries: 1024
Clients: 8, Machine: 4 cores



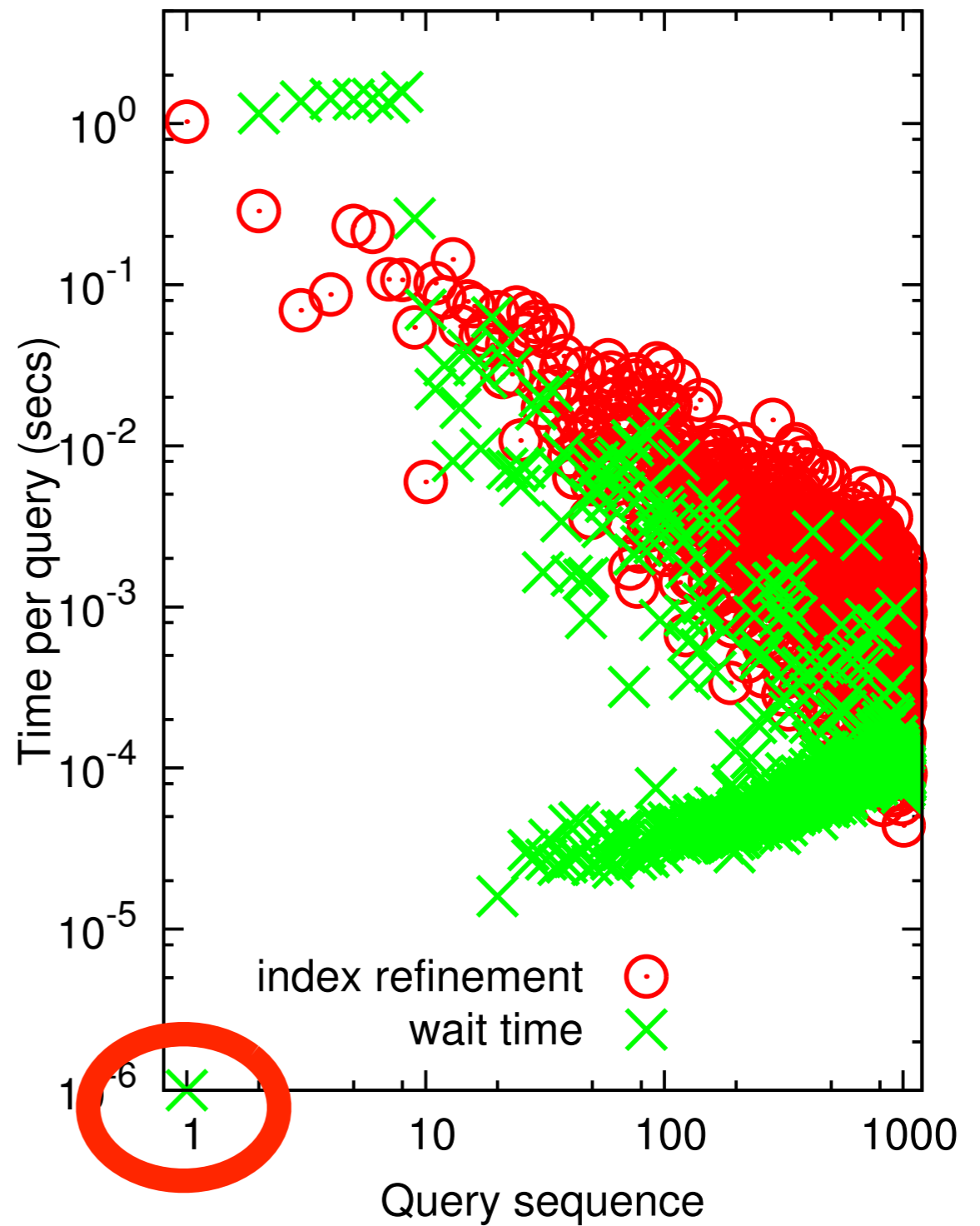
crack time and wait time

Query: select sum(A) from R where v1 < A1 < v2
Selectivity: 0.01%, Random, # of queries: 1024
Clients: 8, Machine: 4 cores



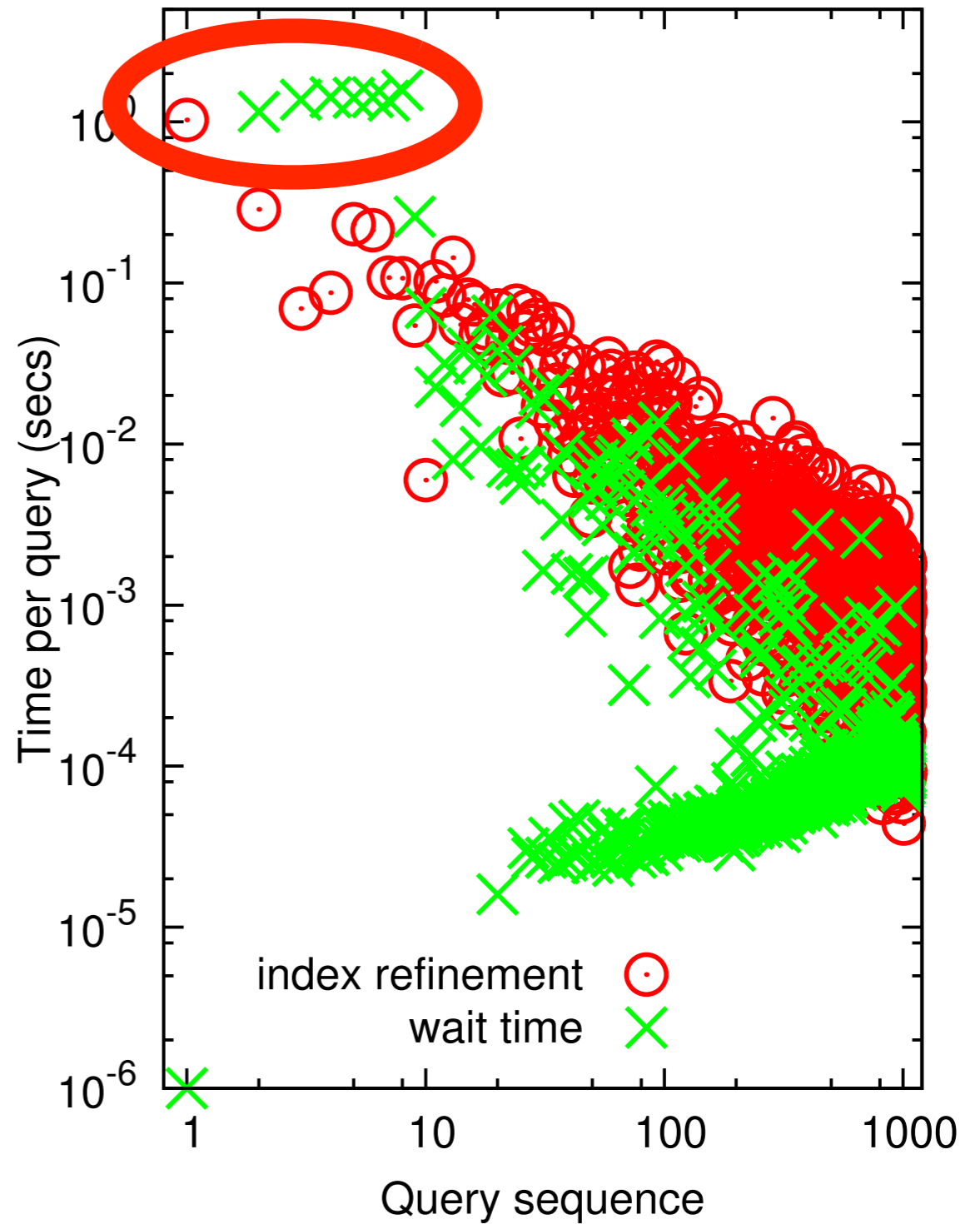
crack time and wait time

Query: select sum(A) from R where v1 < A1 < v2
Selectivity: 0.01%, Random, # of queries: 1024
Clients: 8, Machine: 4 cores



crack time and wait time

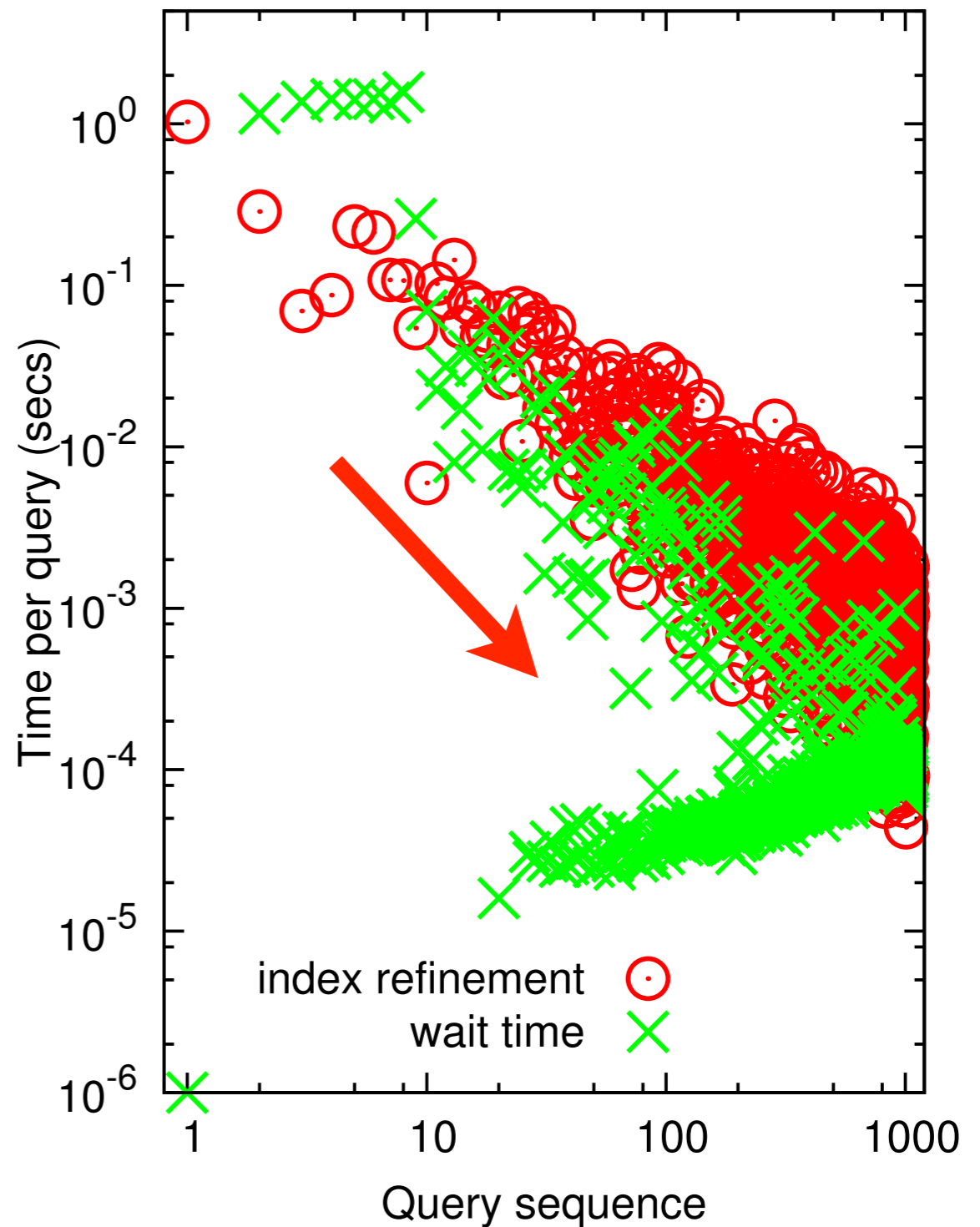
Query: select sum(A) from R where v1 < A1 < v2
Selectivity: 0.01%, Random, # of queries: 1024
Clients: 8, Machine: 4 cores



crack time and wait time

Query: select sum(A) from R where v1 < A1 < v2
Selectivity: 0.01%, Random, # of queries: 1024
Clients: 8, Machine: 4 cores

adaptive indexing
maintains
the adaptive behavior

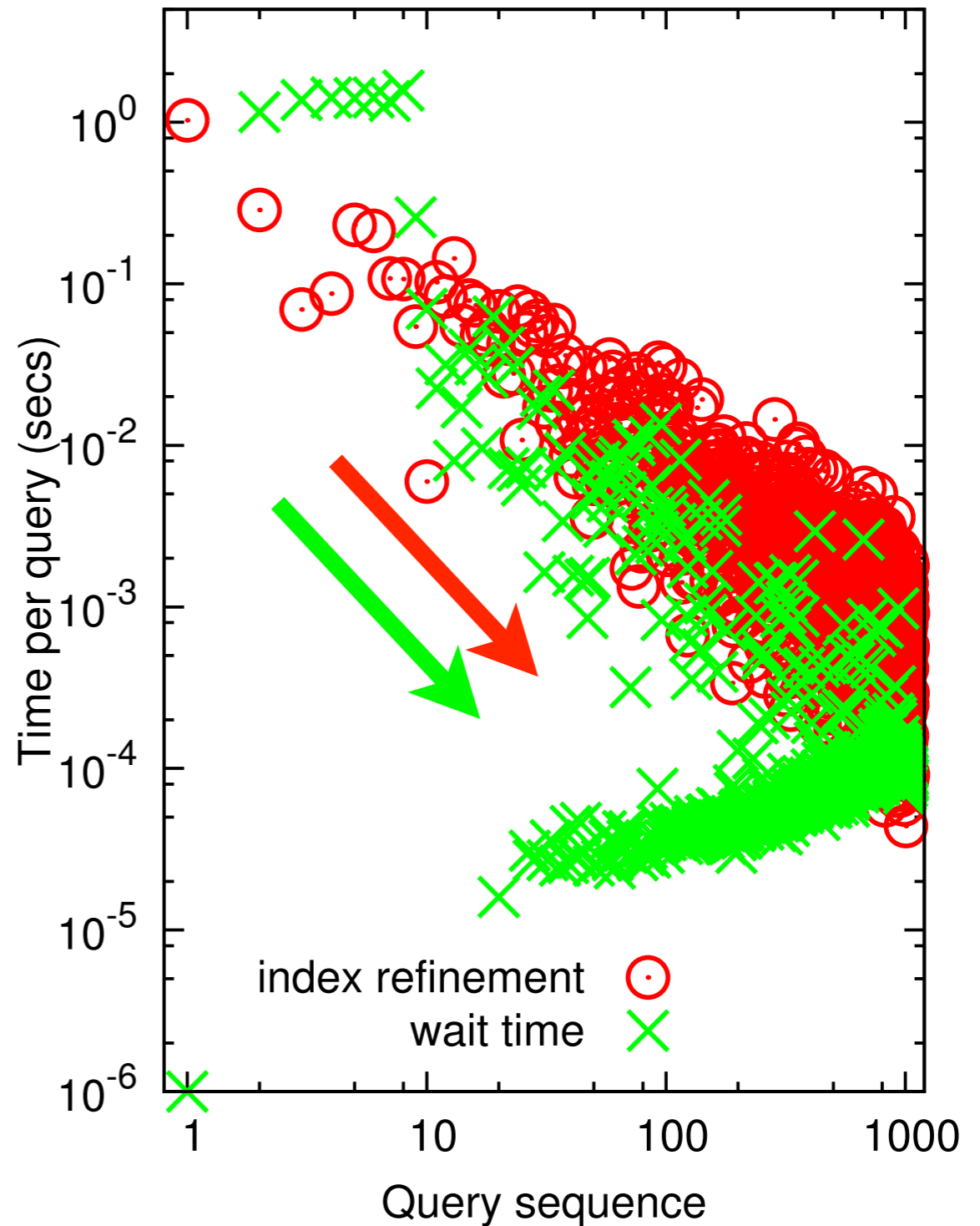


crack time and wait time

Query: select sum(A) from R where v1 < A1 < v2
Selectivity: 0.01%, Random, # of queries: 1024
Clients: 8, Machine: 4 cores

adaptive indexing
maintains
the adaptive behavior

adaptive behavior also for
conflicts

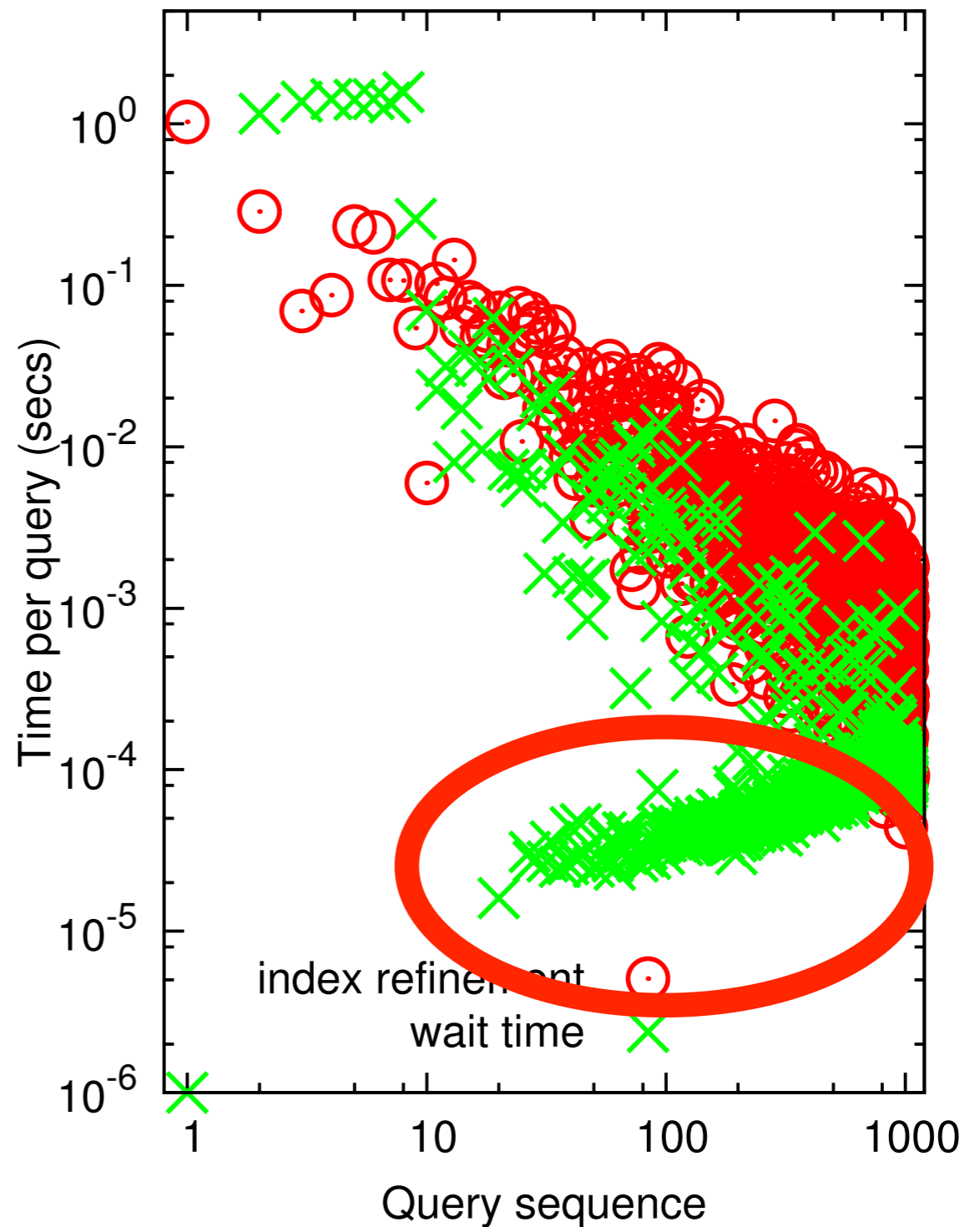


crack time and wait time

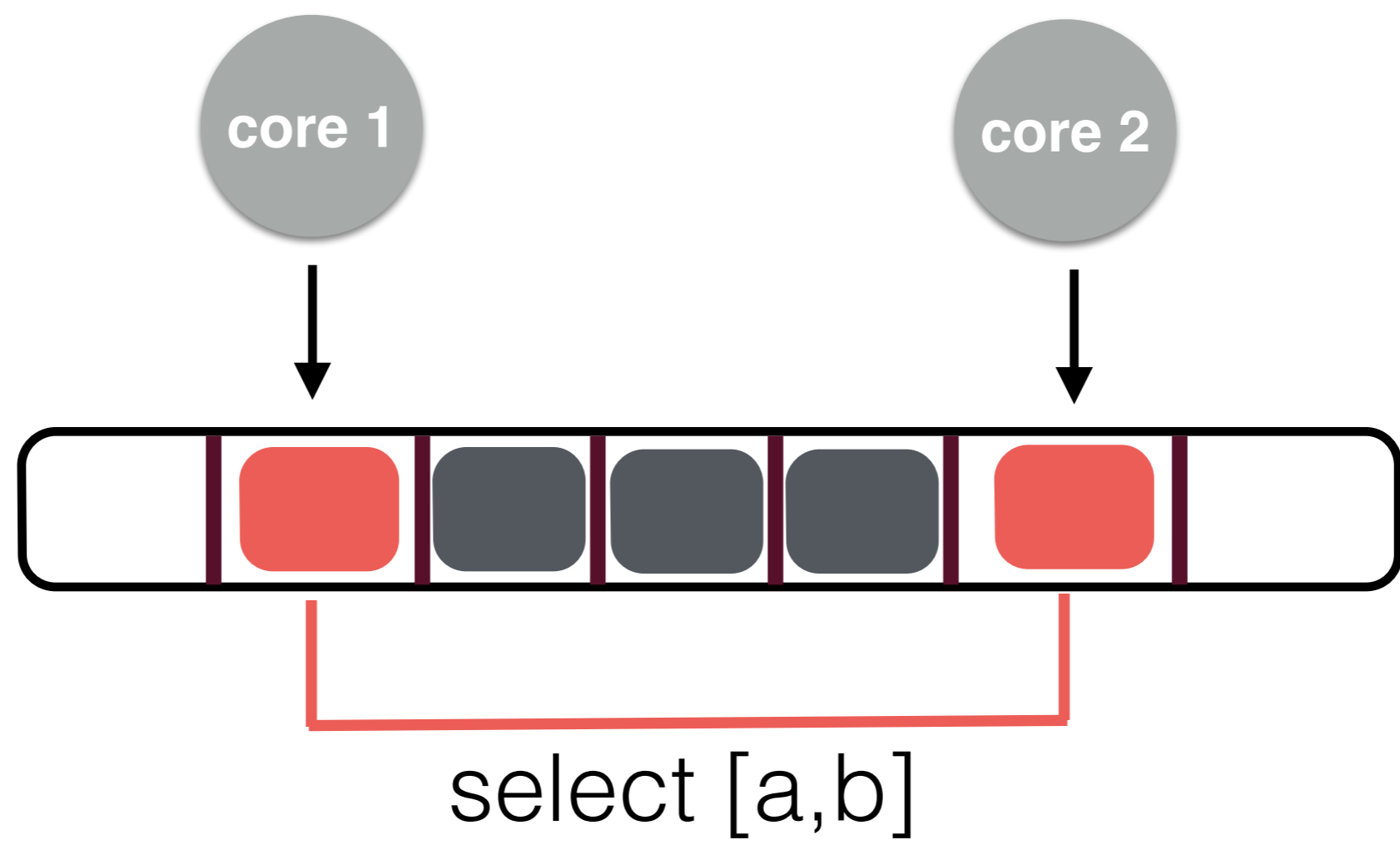
Query: select sum(A) from R where v1 < A1 < v2
Selectivity: 0.01%, Random, # of queries: 1024
Clients: 8, Machine: 4 cores

adaptive indexing
maintains
the adaptive behavior

adaptive behavior also for
conflicts



multi-core utilization





select [a,b]



select [a,b]





select [a,b]





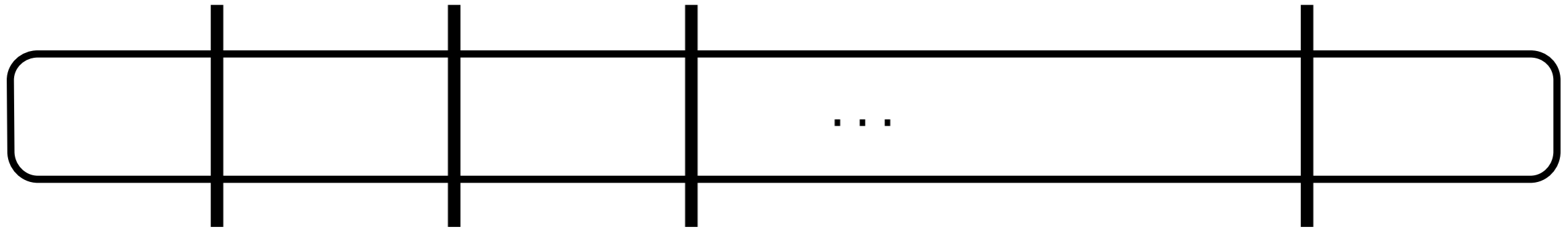
select [a,b]

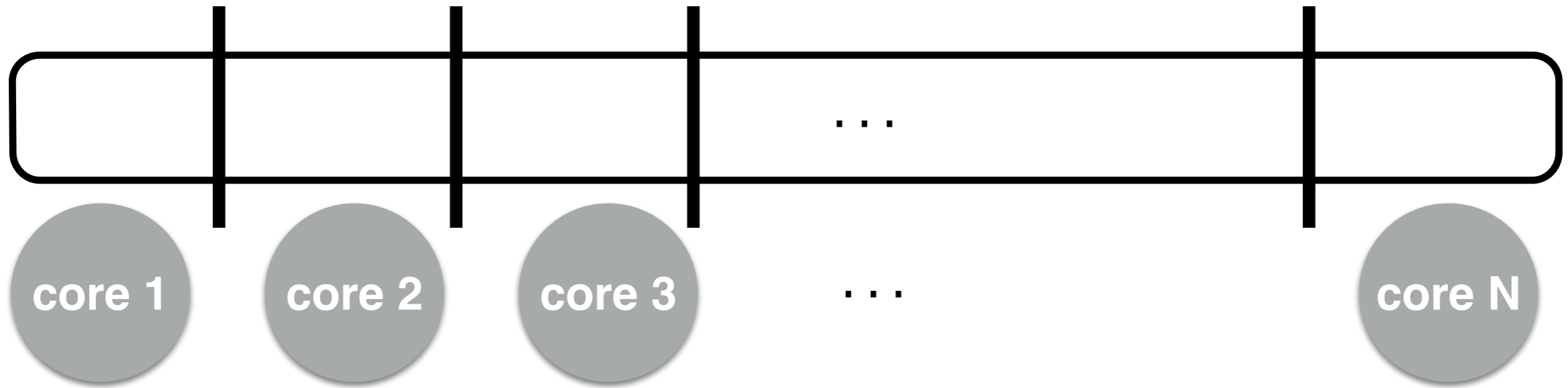


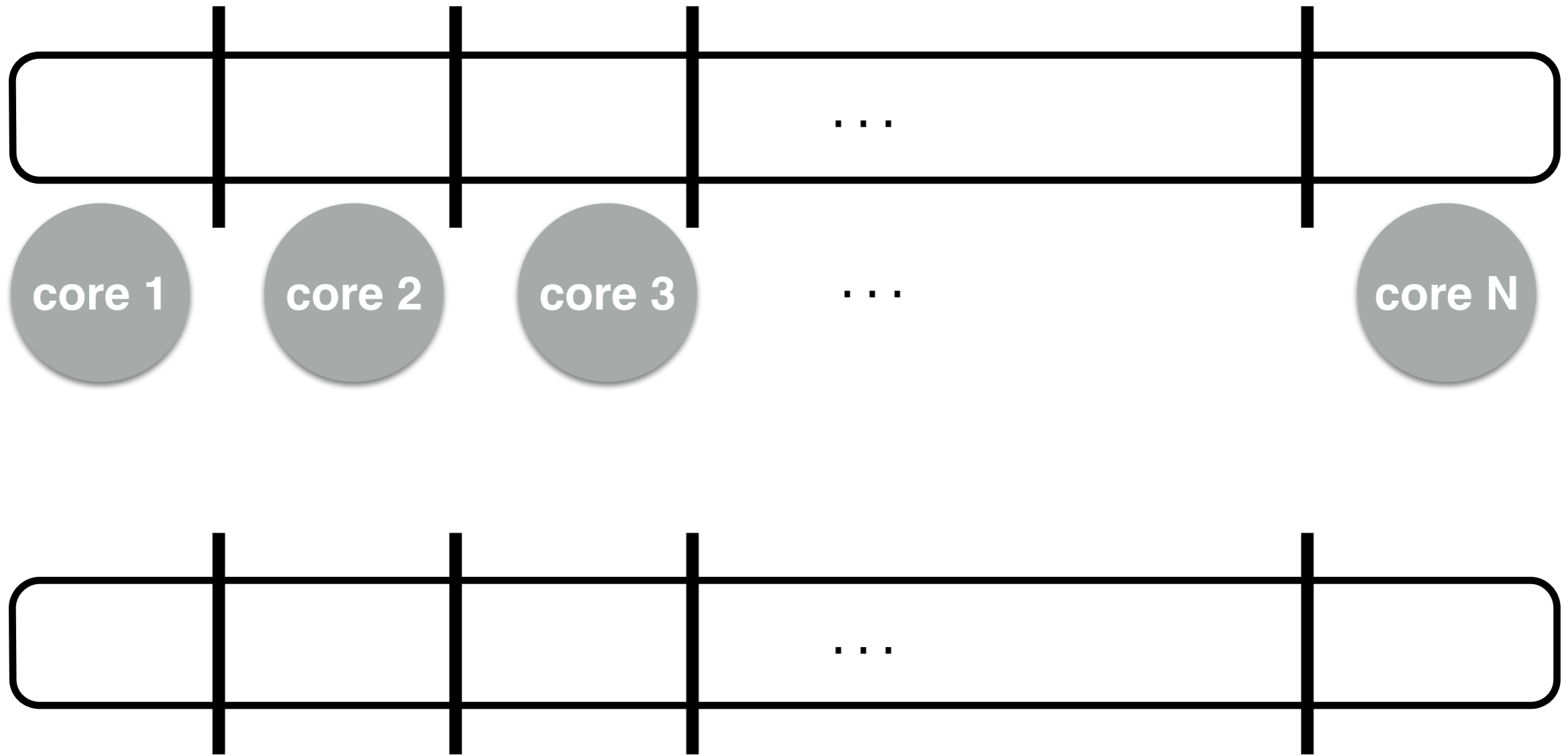
...

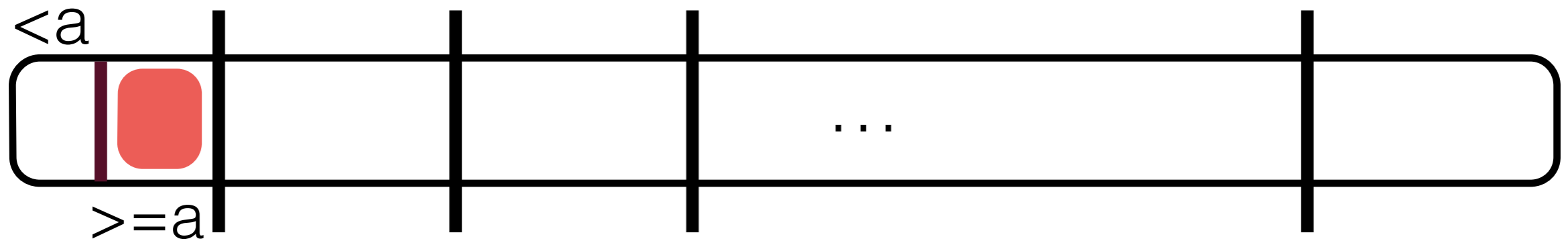
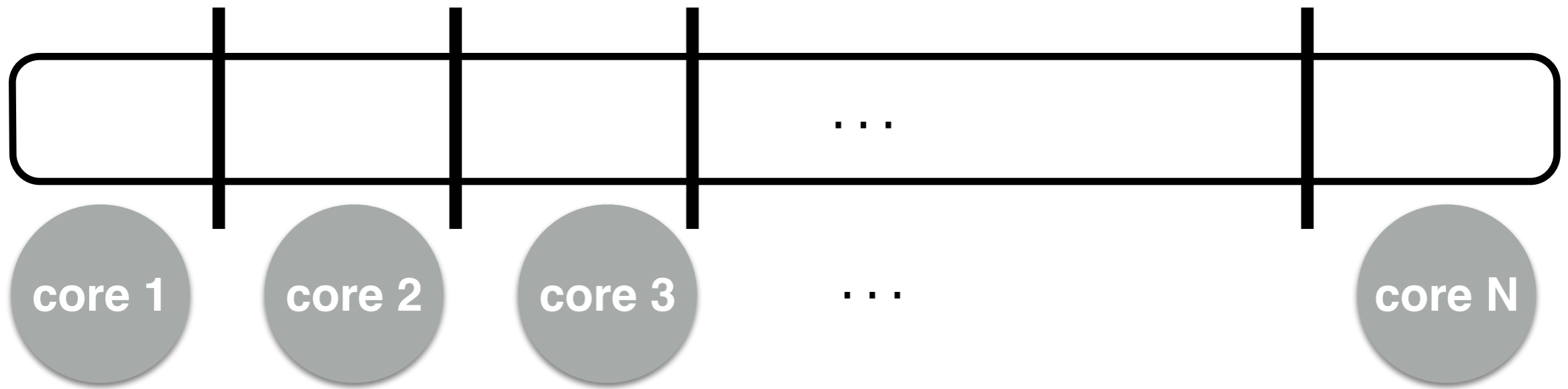


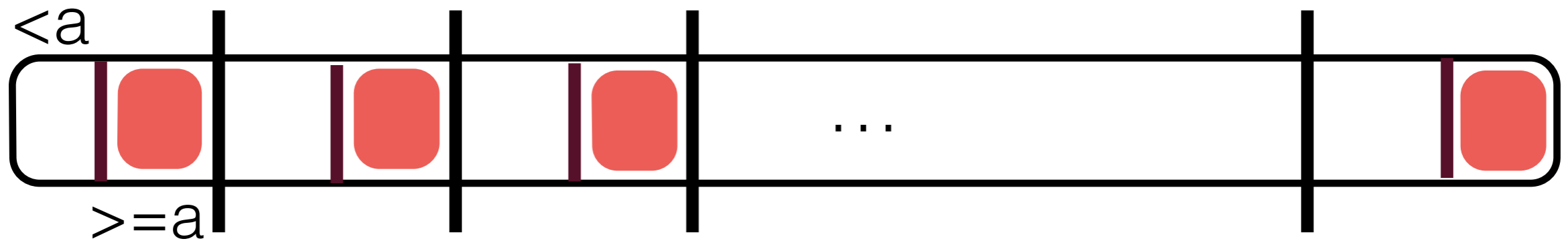
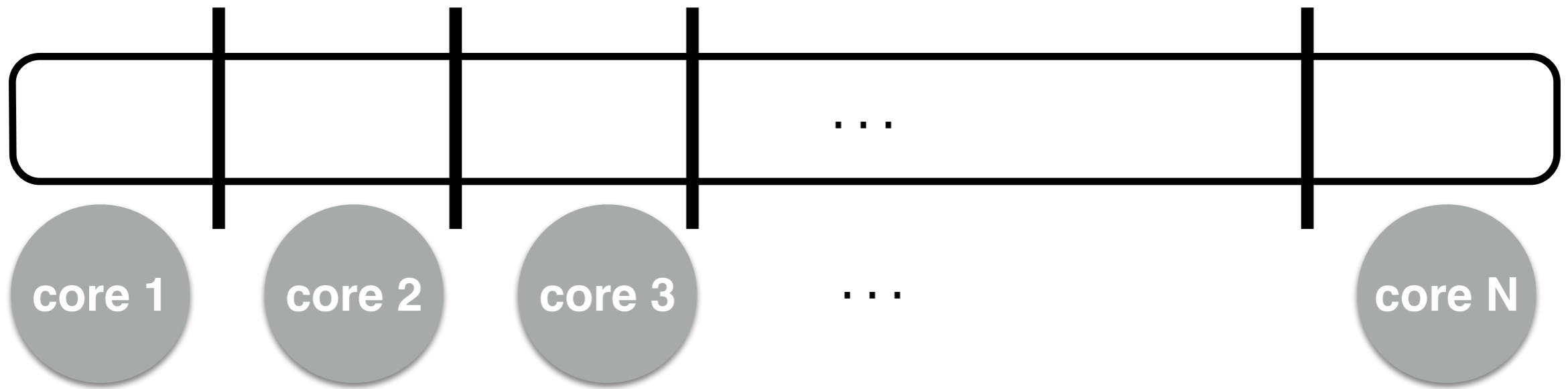


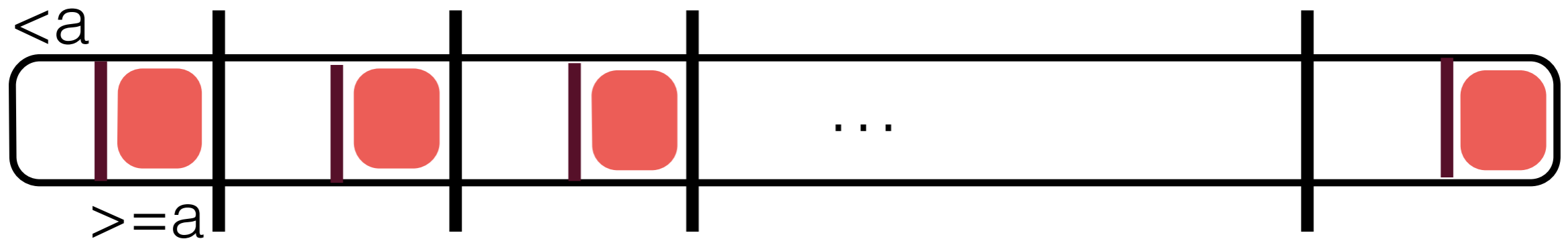
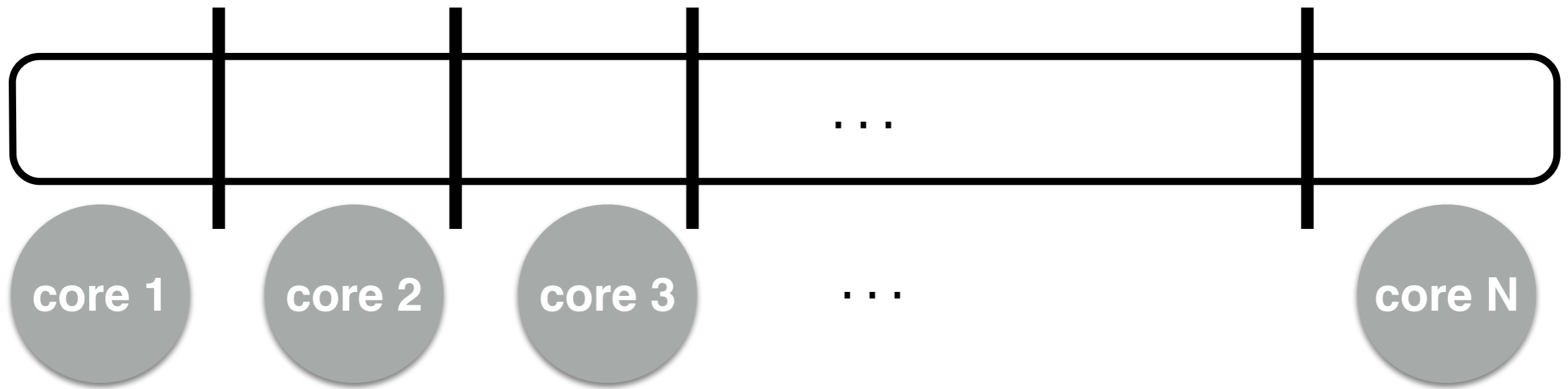




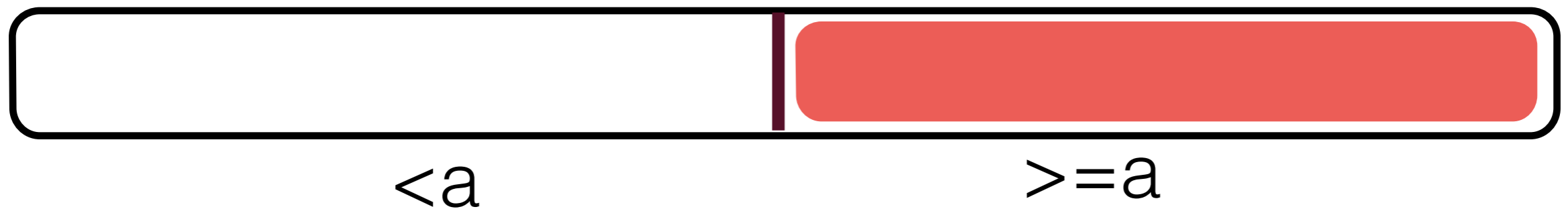









 merge

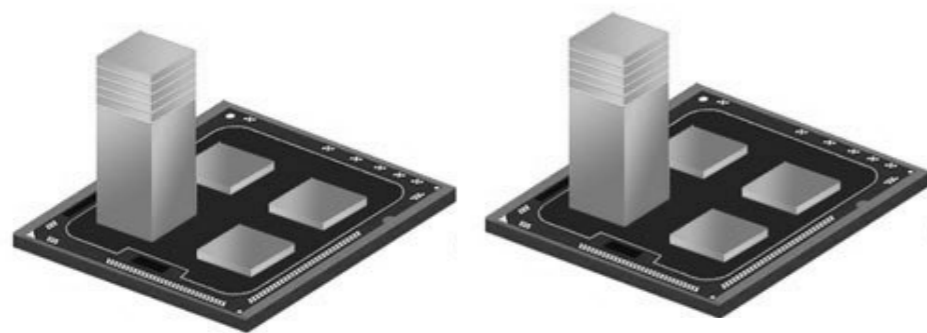


holistic indexing

problem: cores may be under utilized

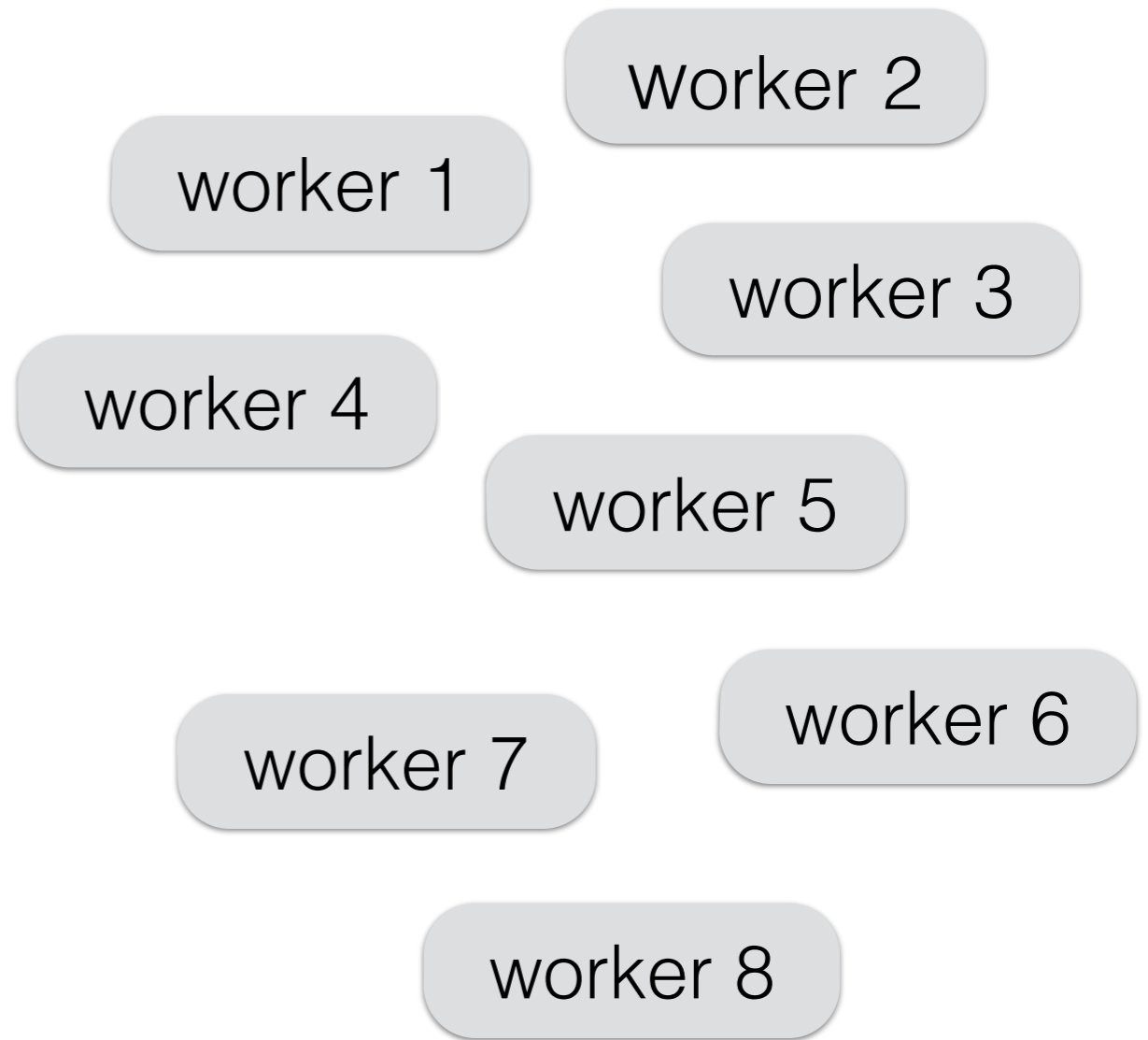
goal: either fully utilize a core or shut it down

scheduler

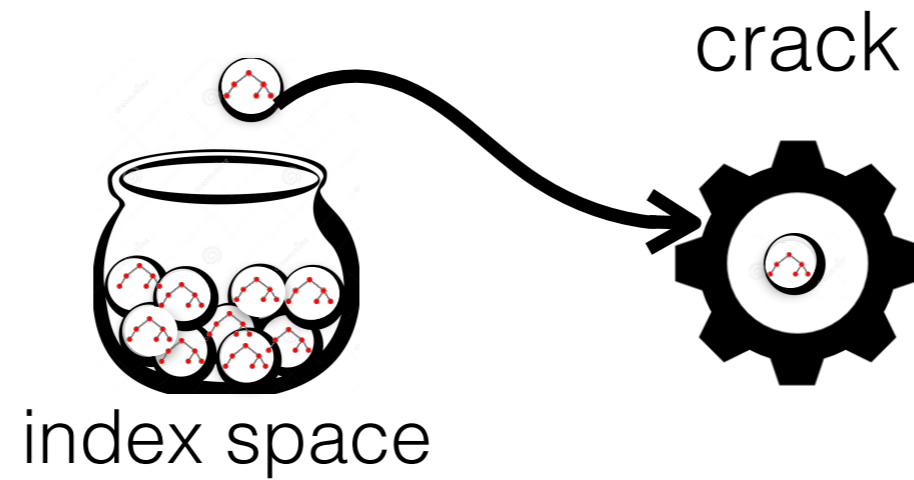


monitoring CPU utilization

thread pool



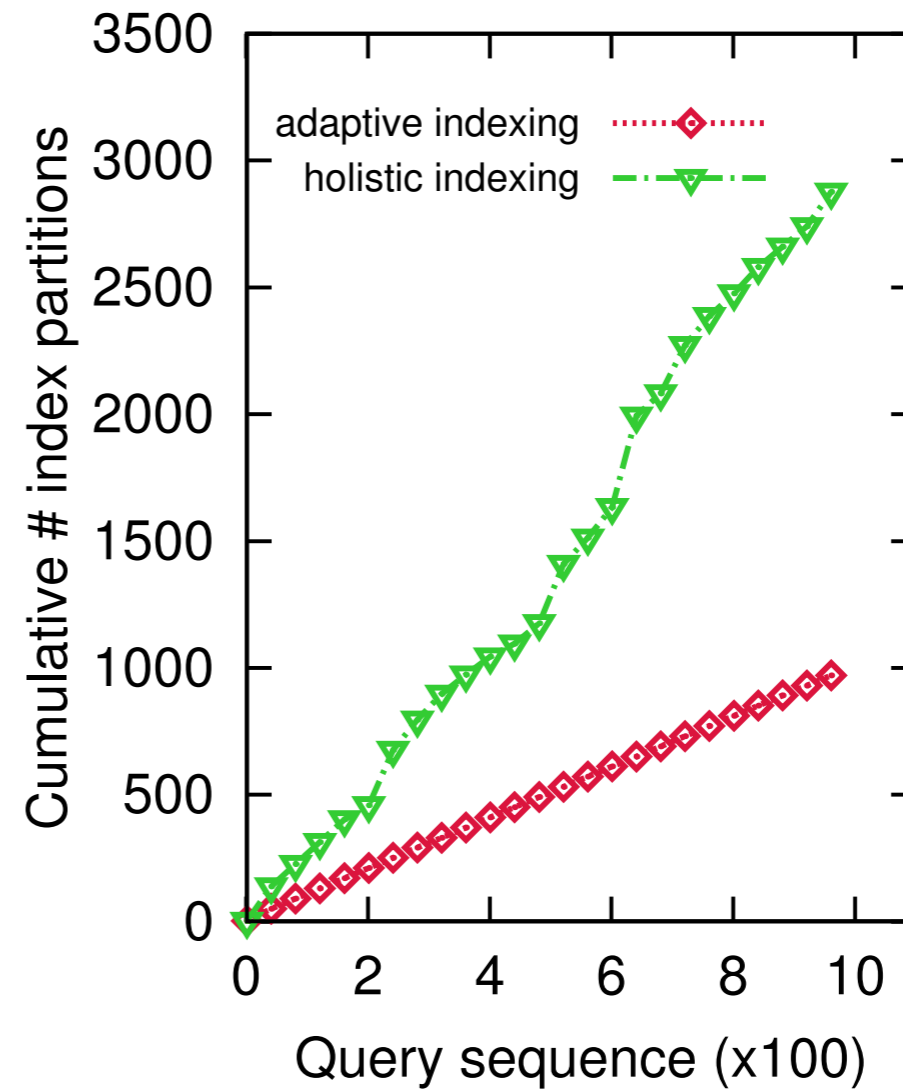
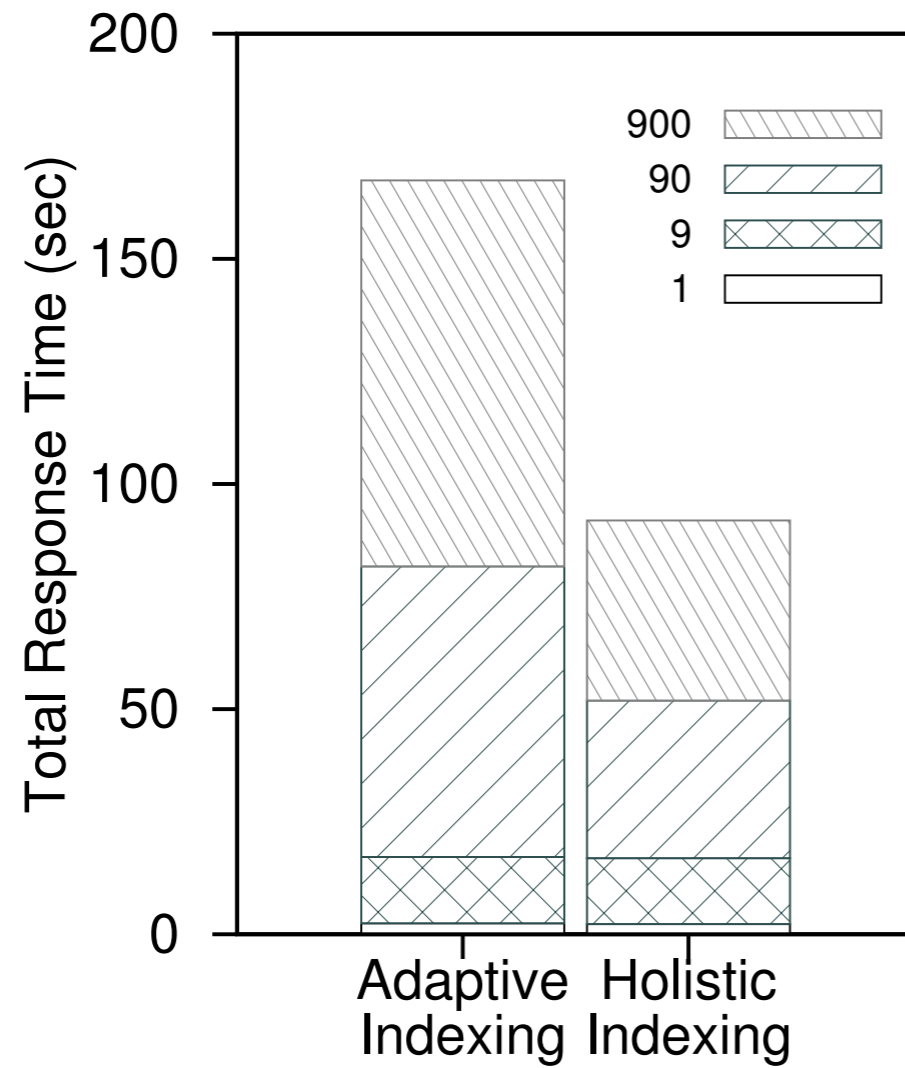
when there is an underutilized CPU, pin a thread to it and to do a cracking task



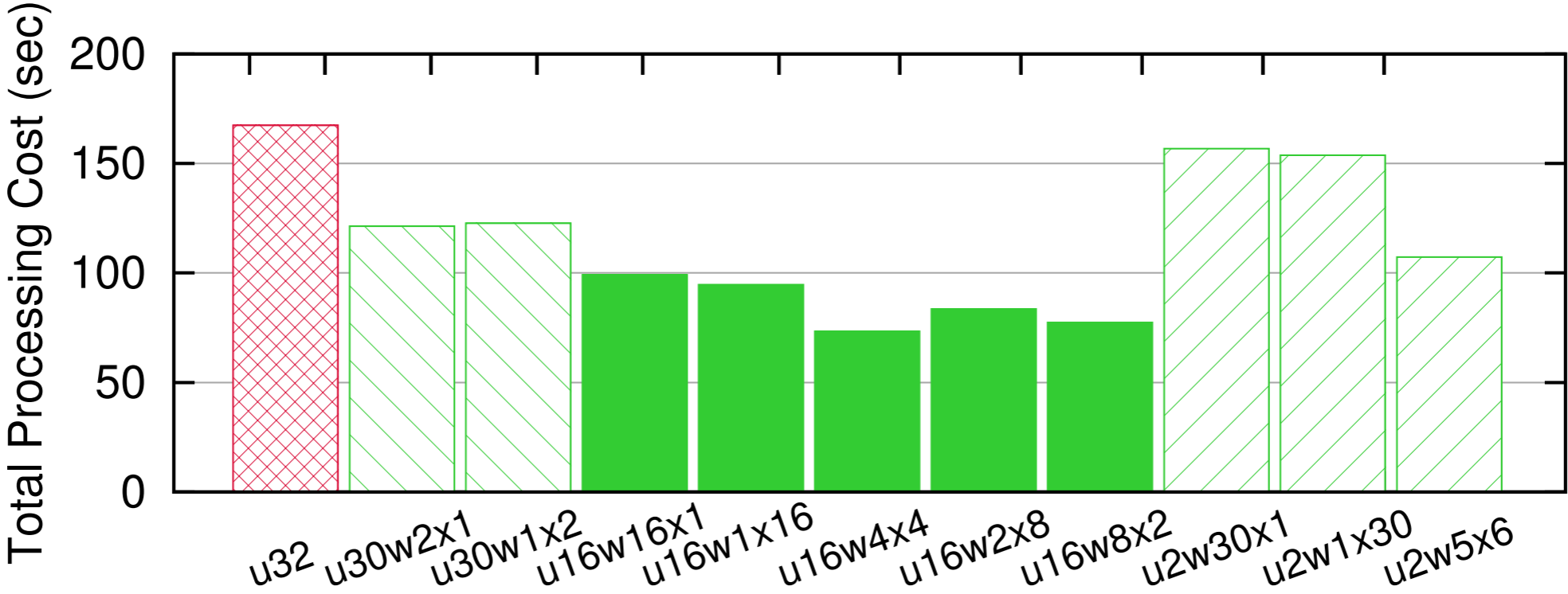
partitions size - access frequency - hit ratio

random works best

10⁸ tuples - 10 attributes, random queries

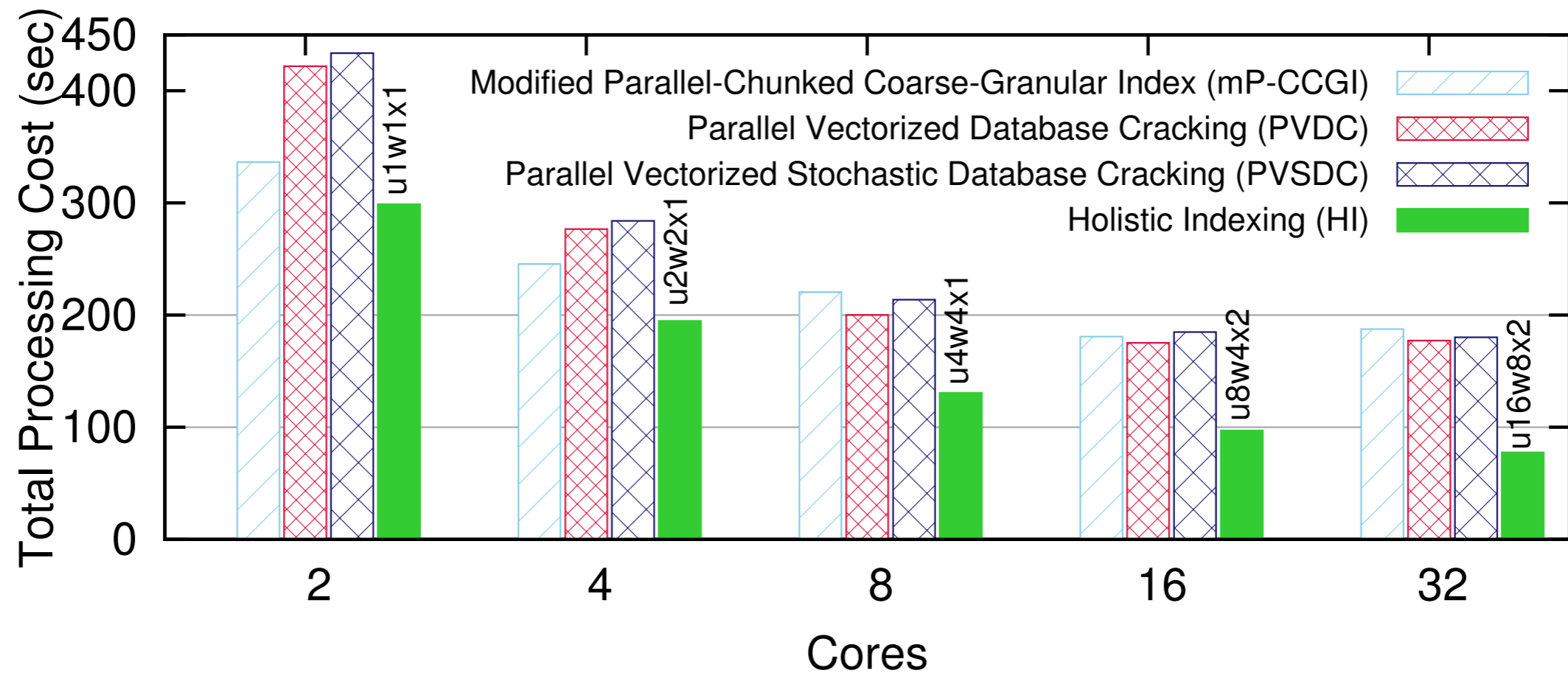


how to distribute the cores - 16 core machine with hyper threading (16+16 hardware threads)



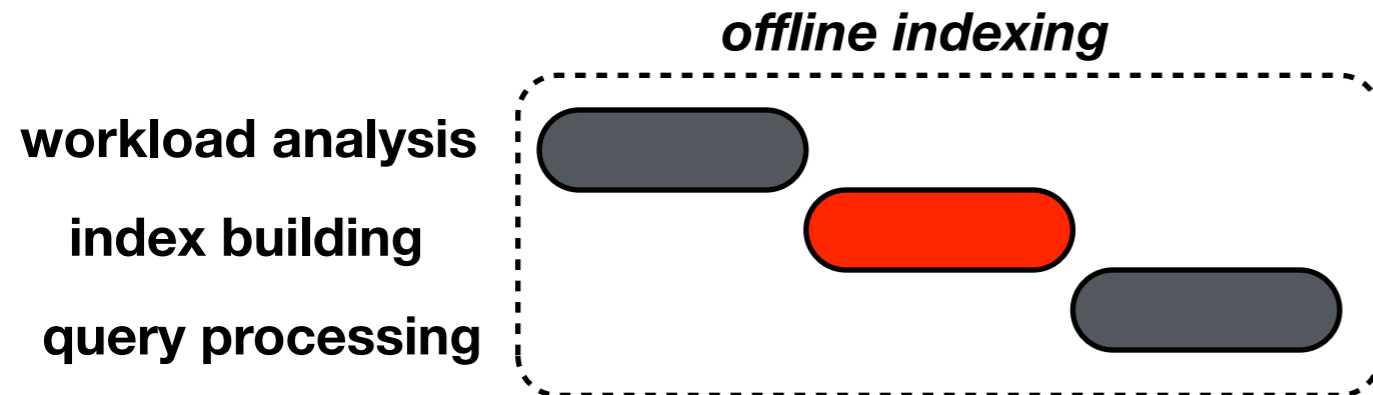
Distribution of 32 Threads between Users and Workers

holistic indexing against using all cores for multi-core adaptive indexing



design space of adaptive indexing

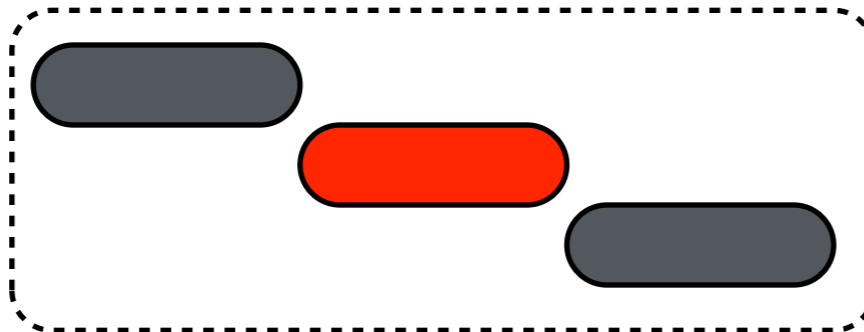
Indexing Overview



Indexing Overview

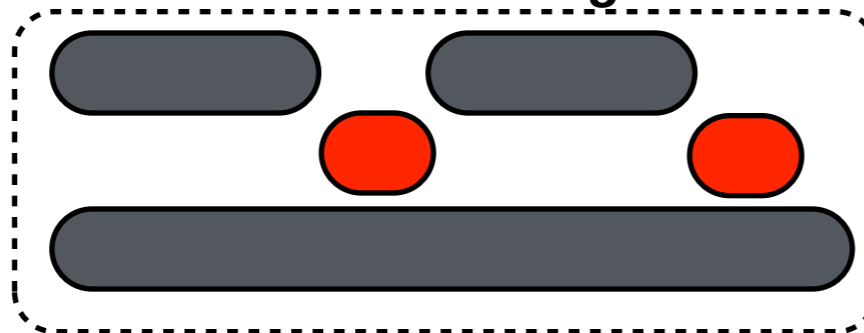
offline indexing

workload analysis
index building
query processing



online indexing

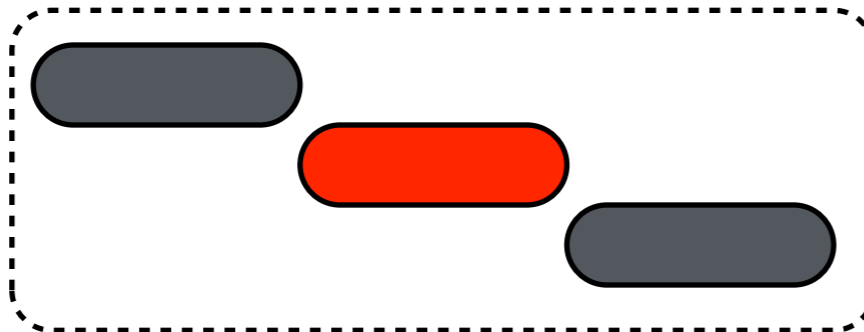
workload analysis
index building
query processing



Indexing Overview

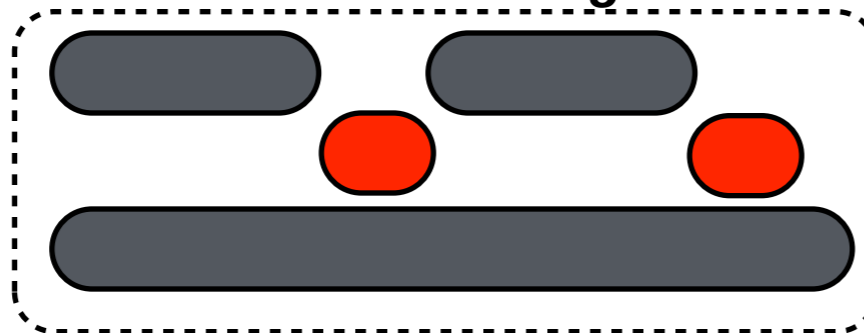
offline indexing

workload analysis
index building
query processing



online indexing

workload analysis
index building
query processing

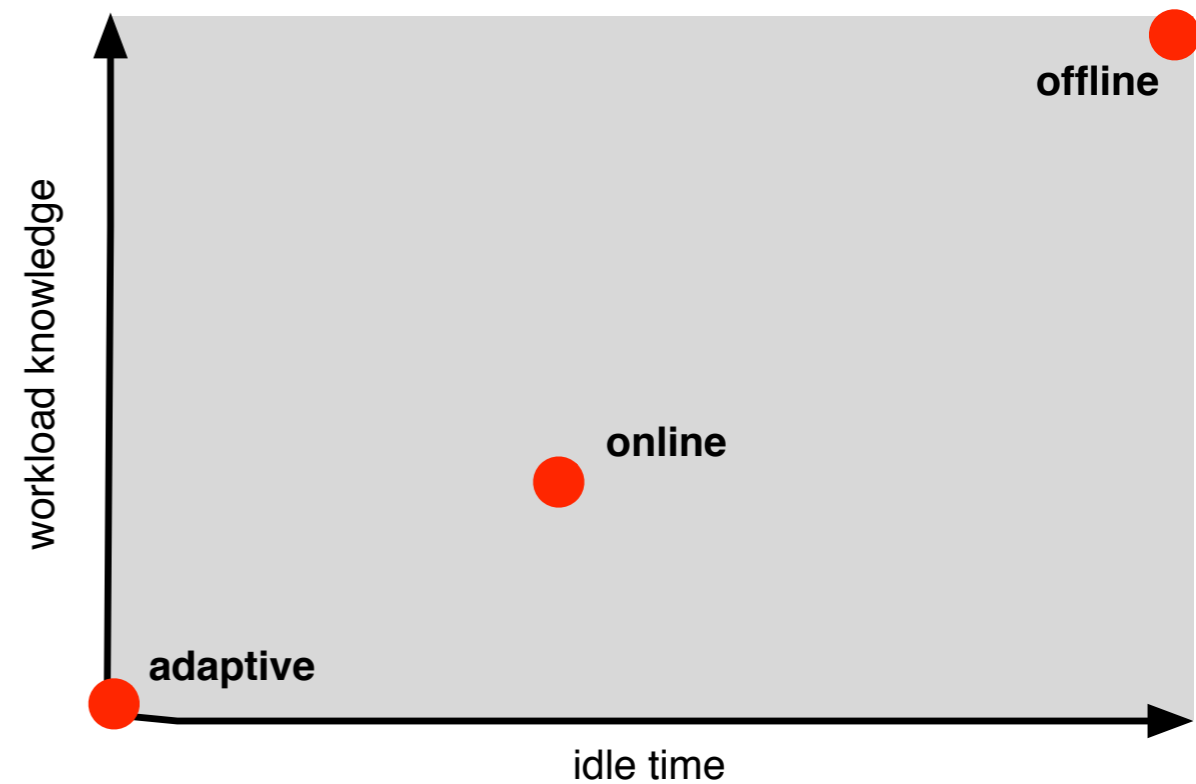
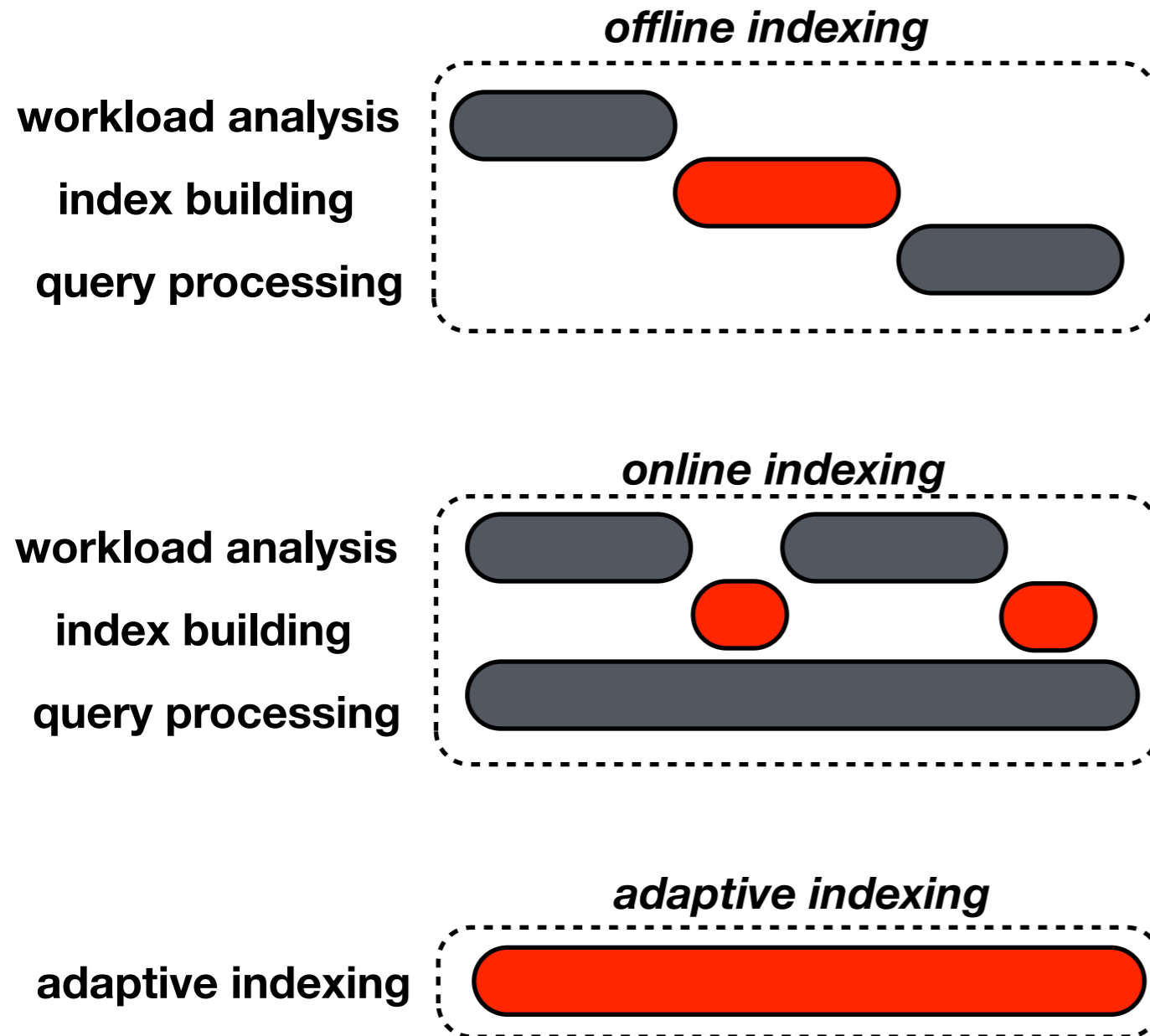


adaptive indexing

adaptive indexing



Indexing Overview



adaptive merging

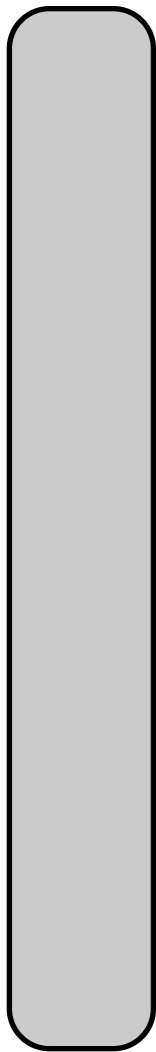
EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

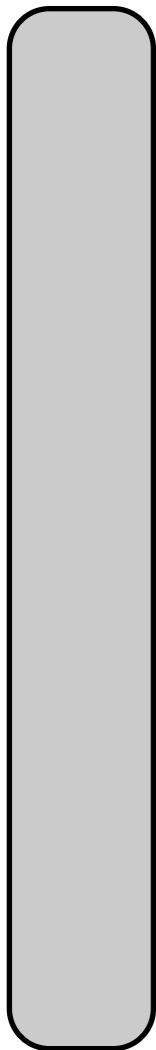


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

```
select(A,50,100)
```

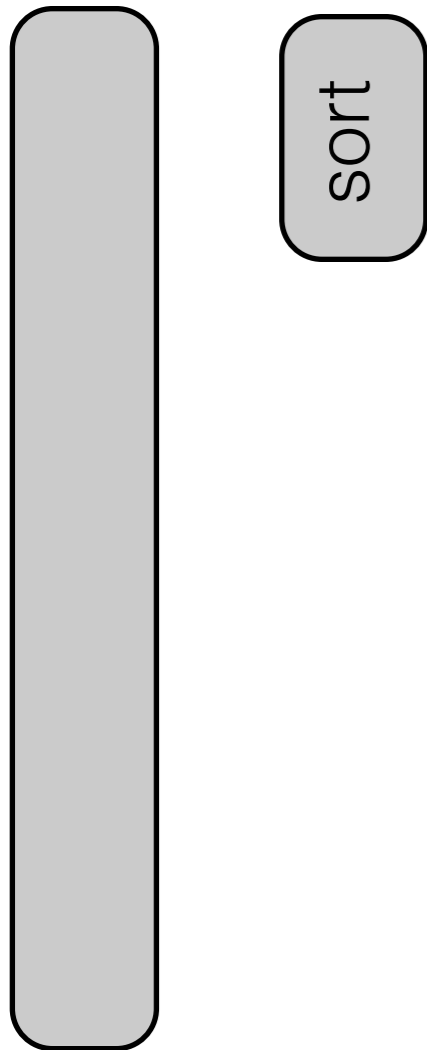


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`

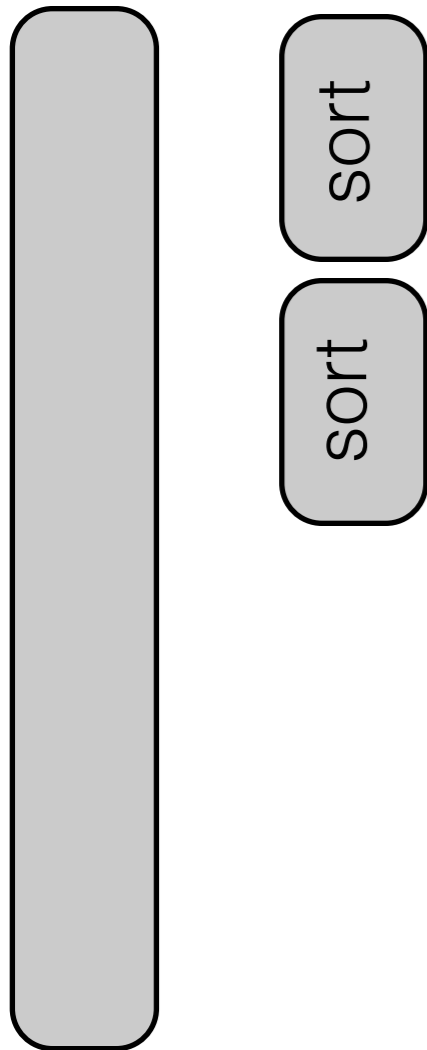


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`

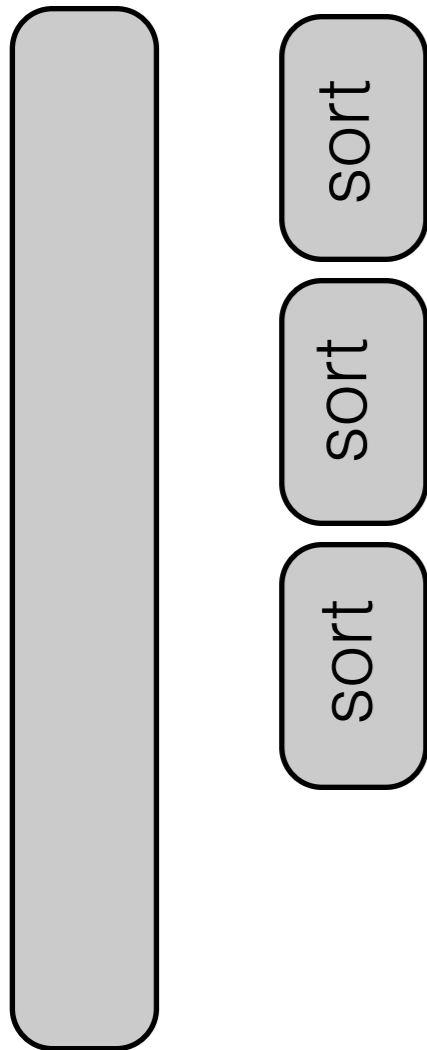


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

```
select(A, 50, 100)
```

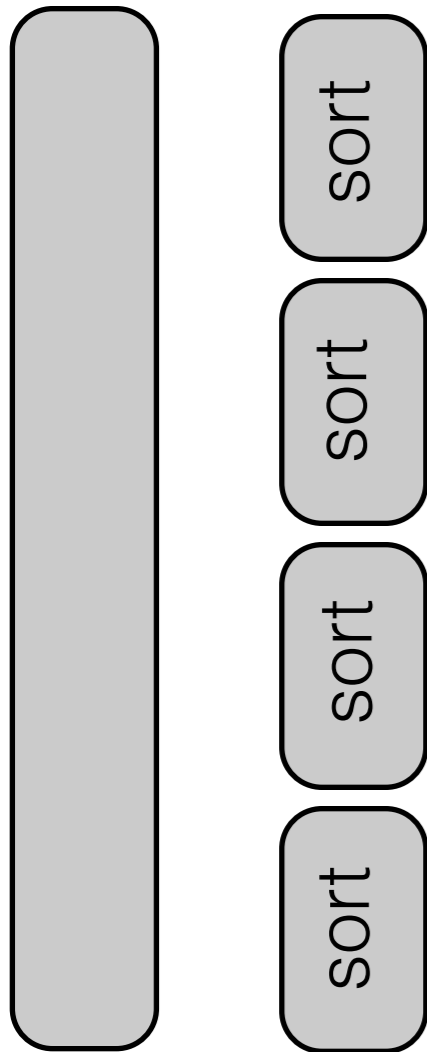


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`

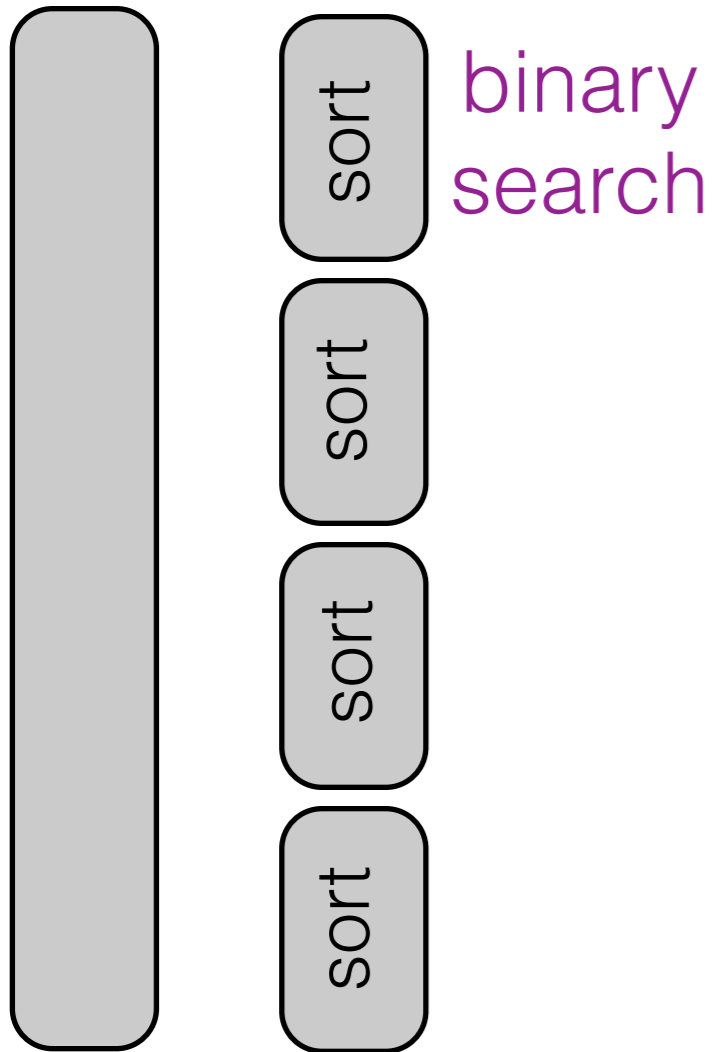


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`

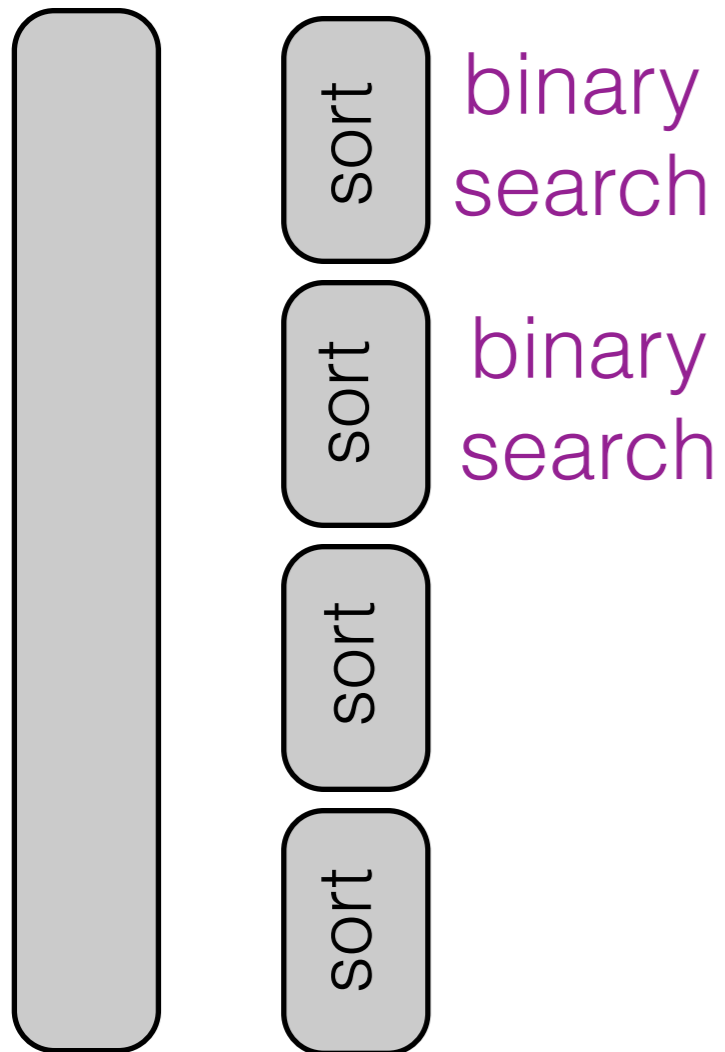


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

select(A,50,100)

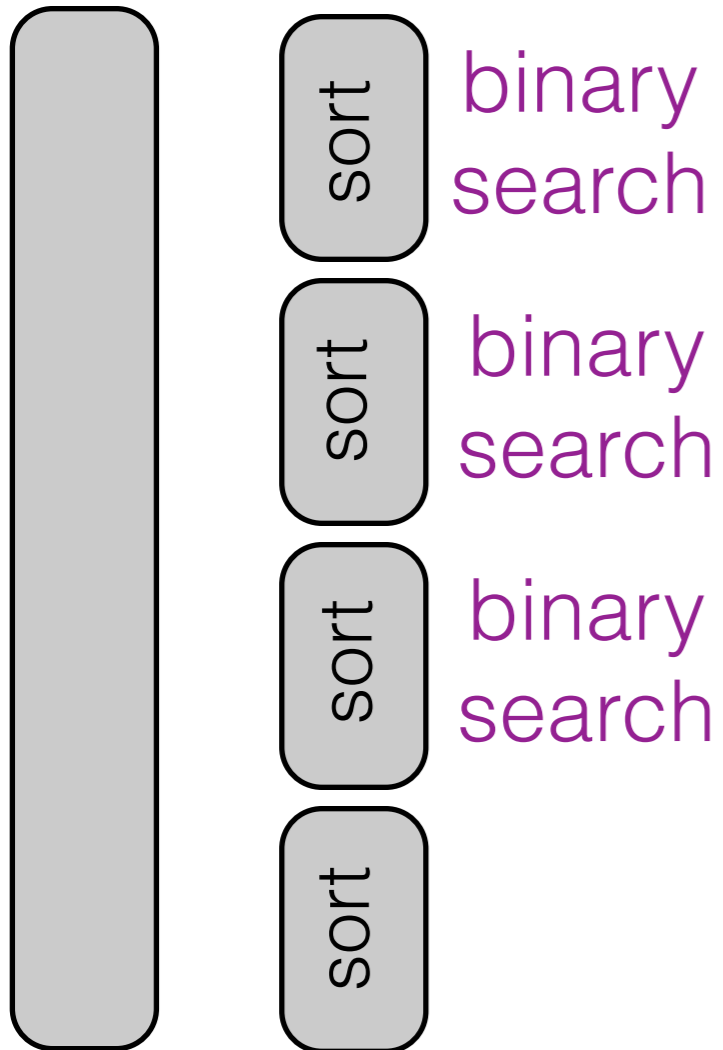


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`

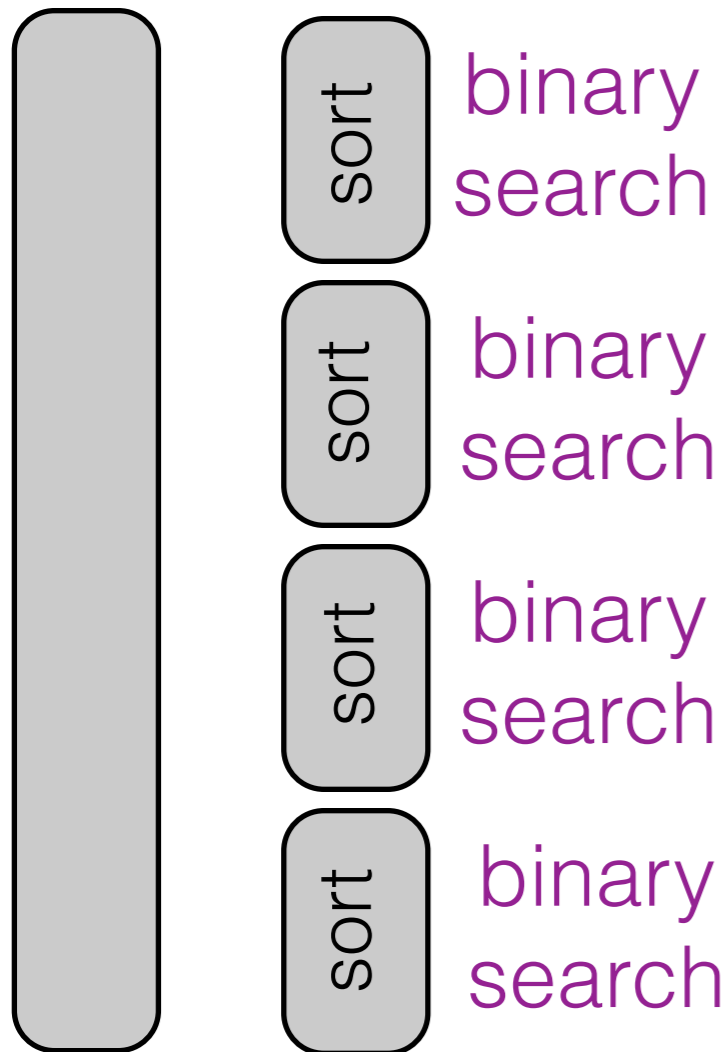


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`

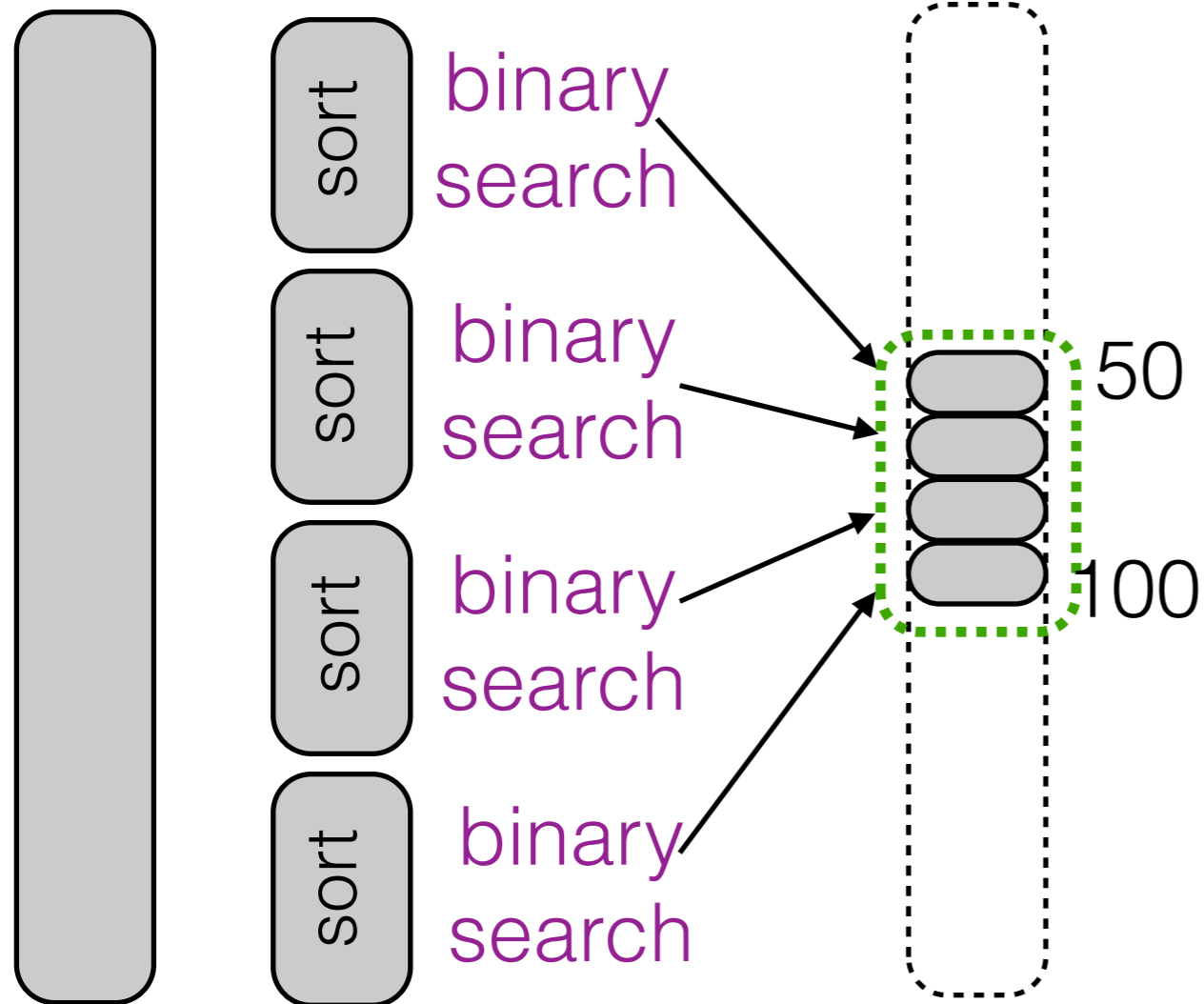


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A, 50, 100)`

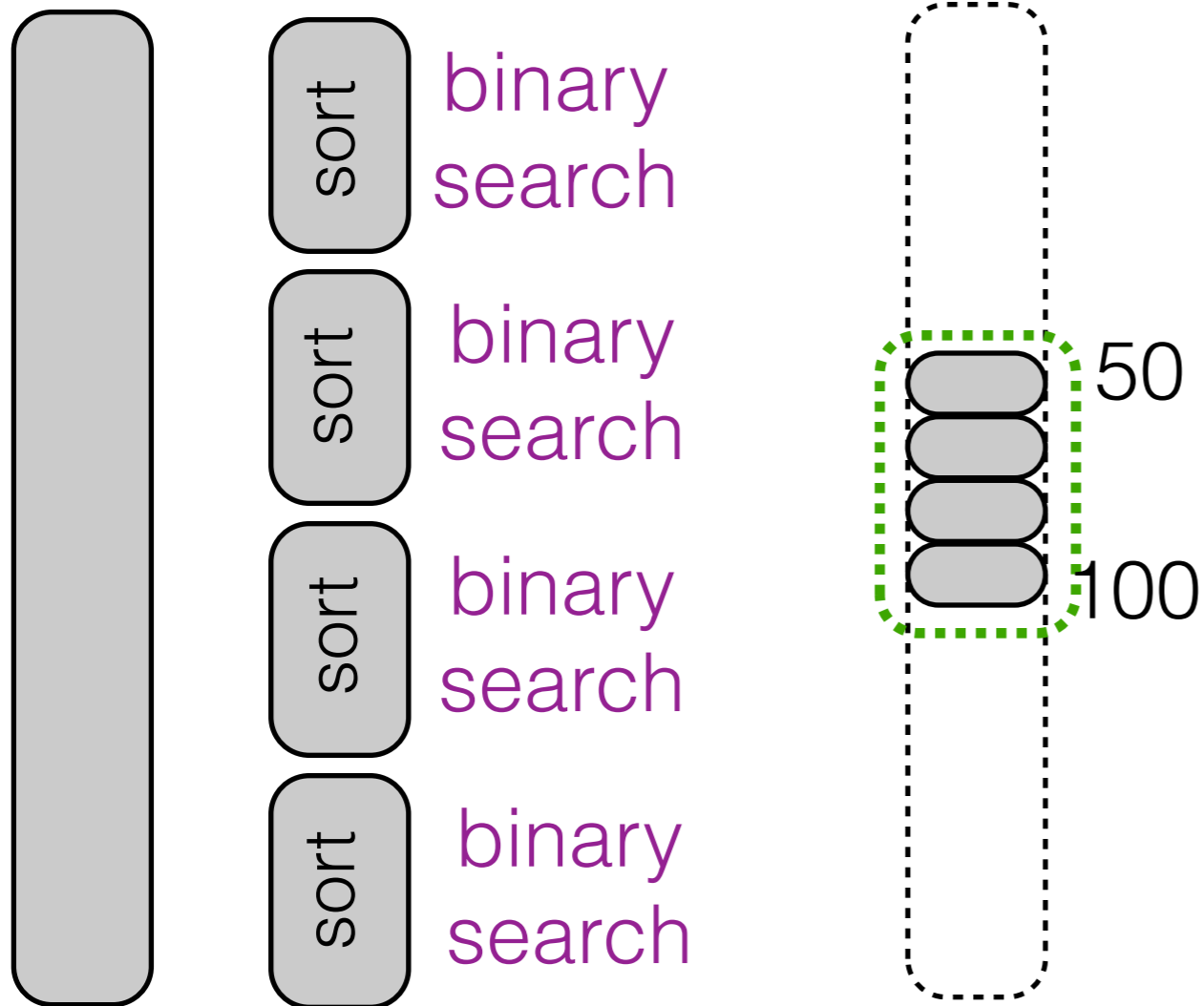


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`

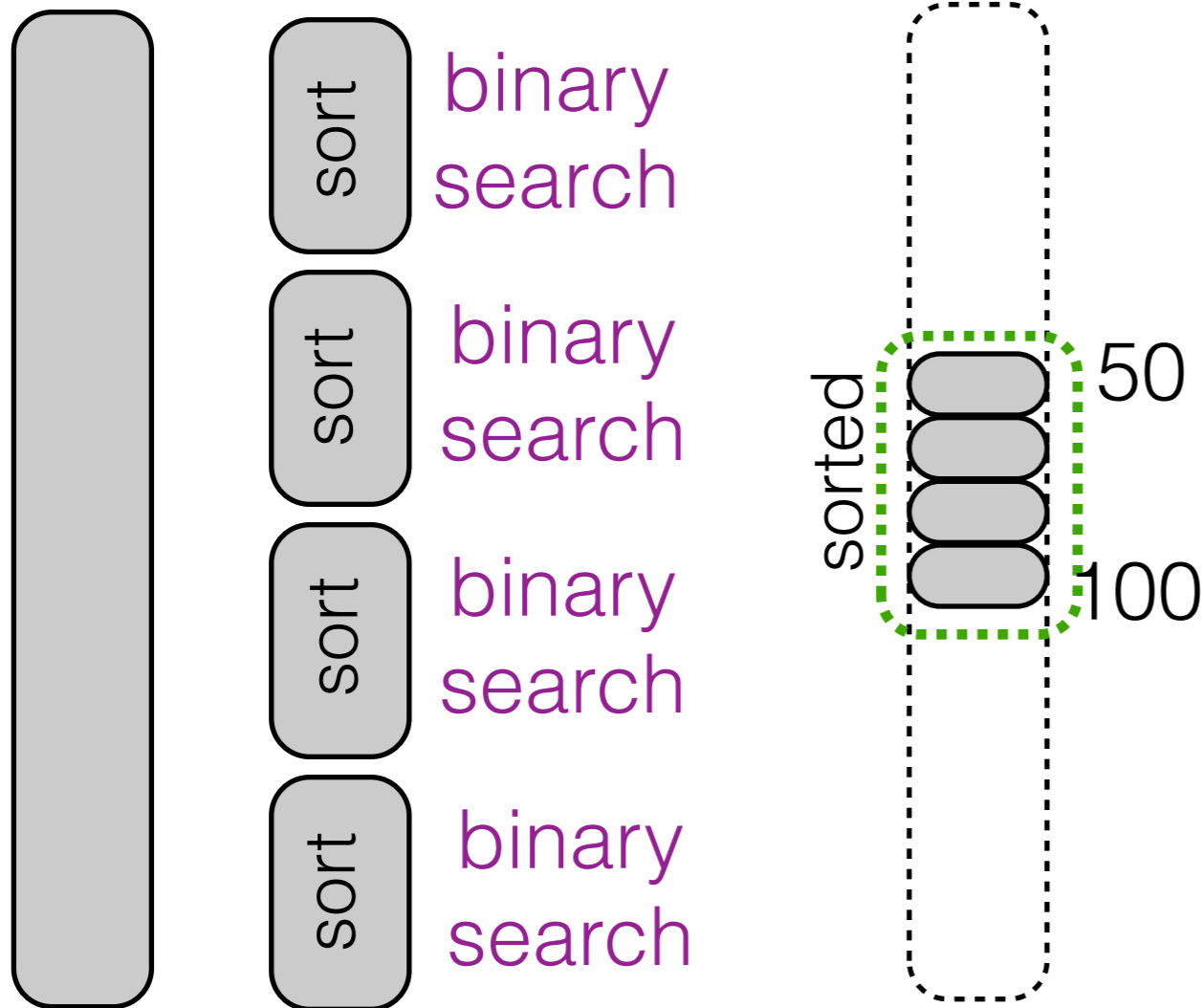


adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`



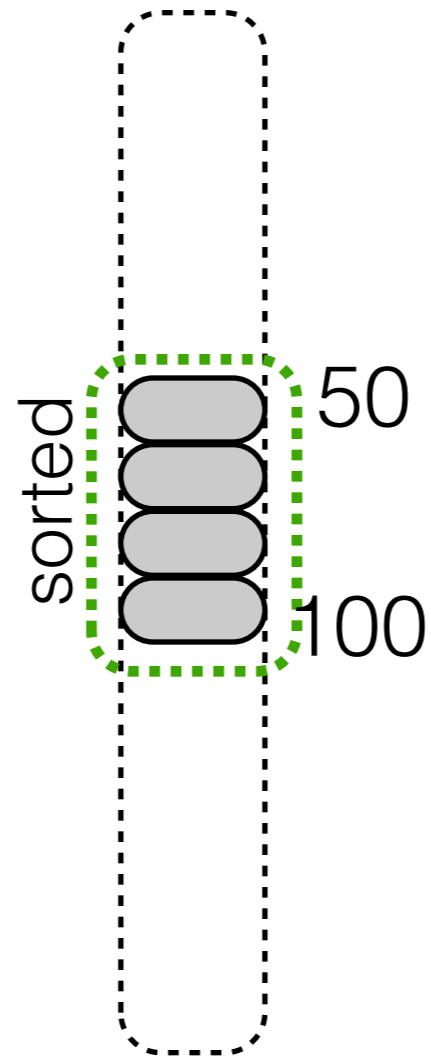
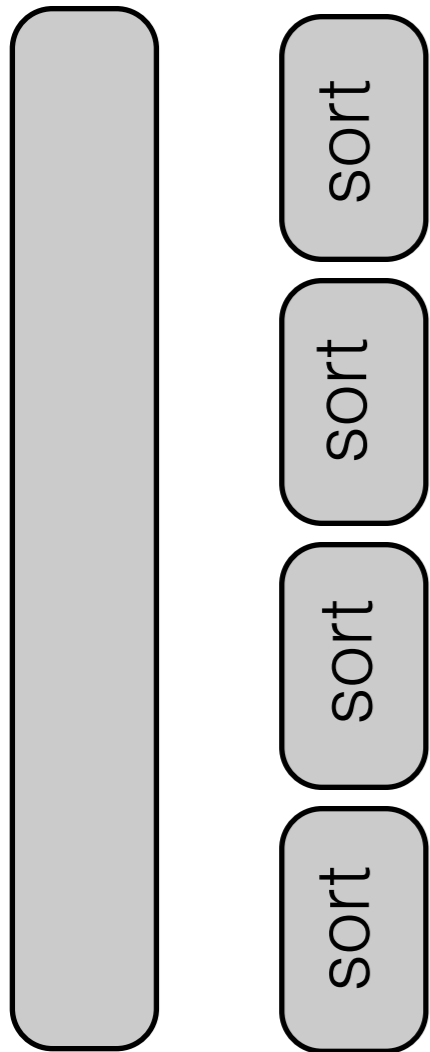
adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`

`select(A,55,70)`



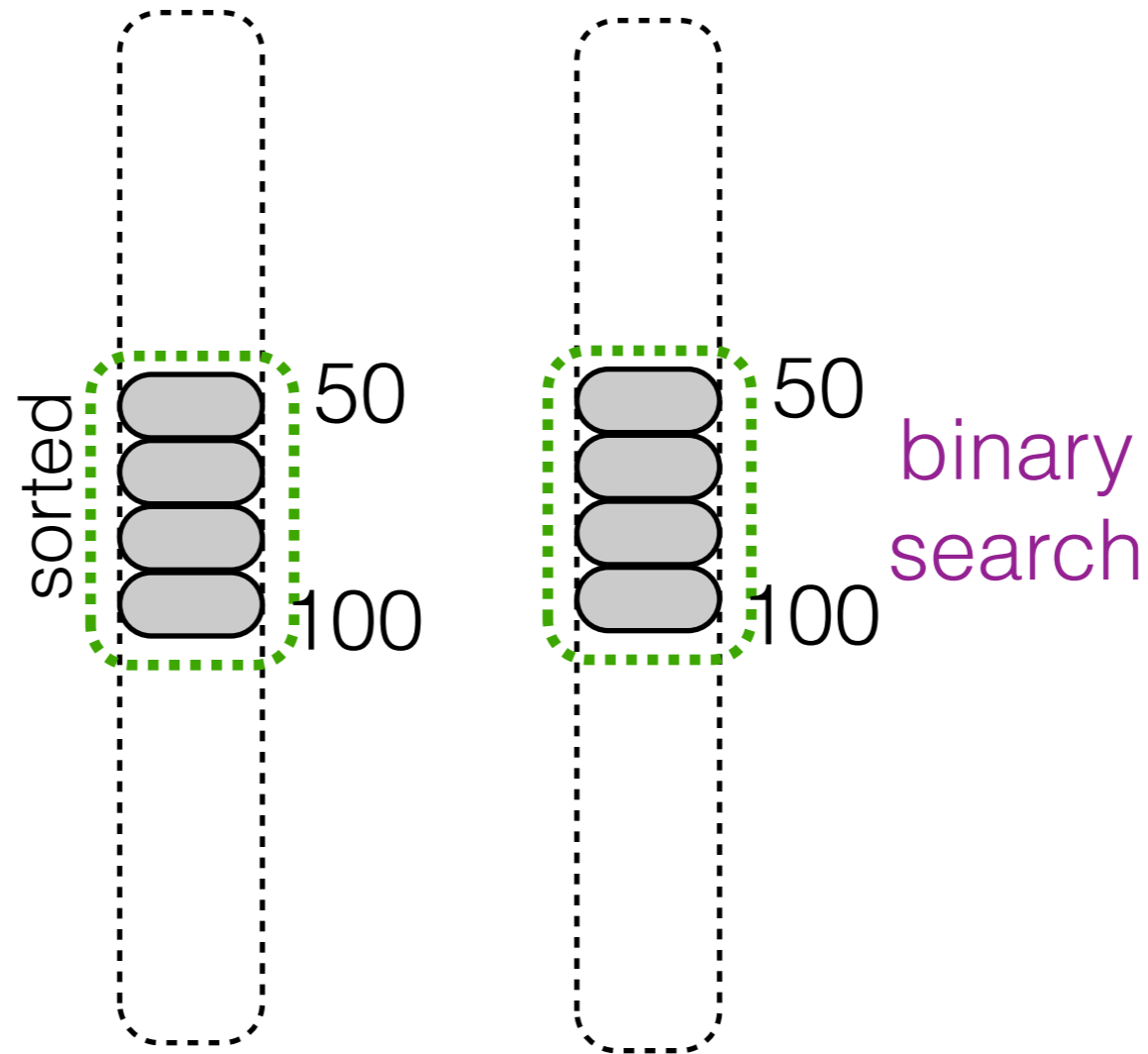
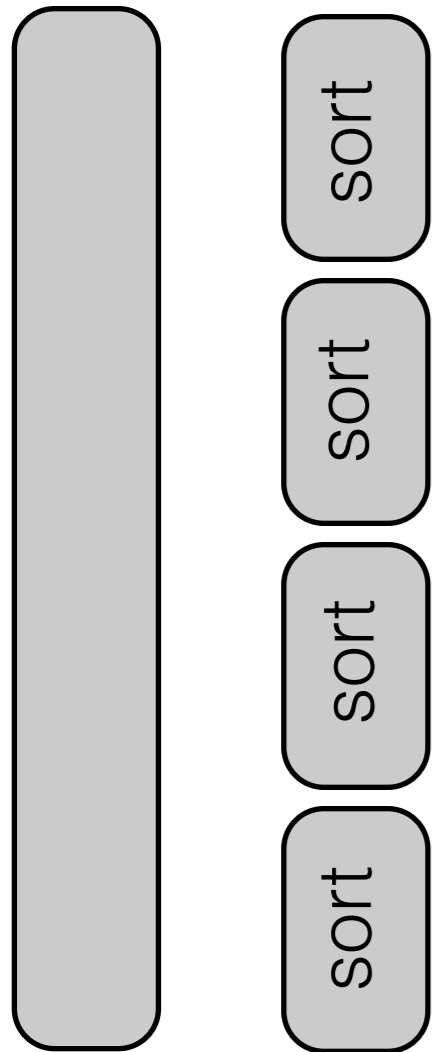
adaptive merging

EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`

`select(A,55,70)`



adaptive merging

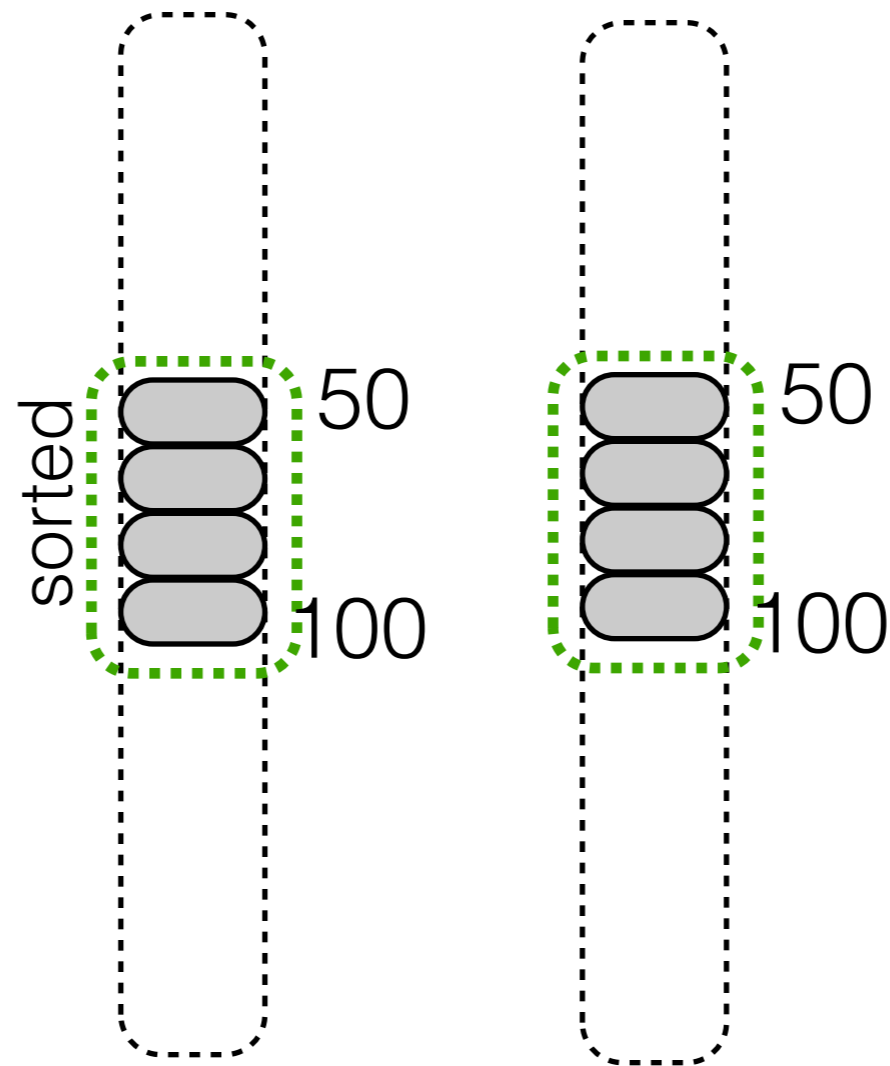
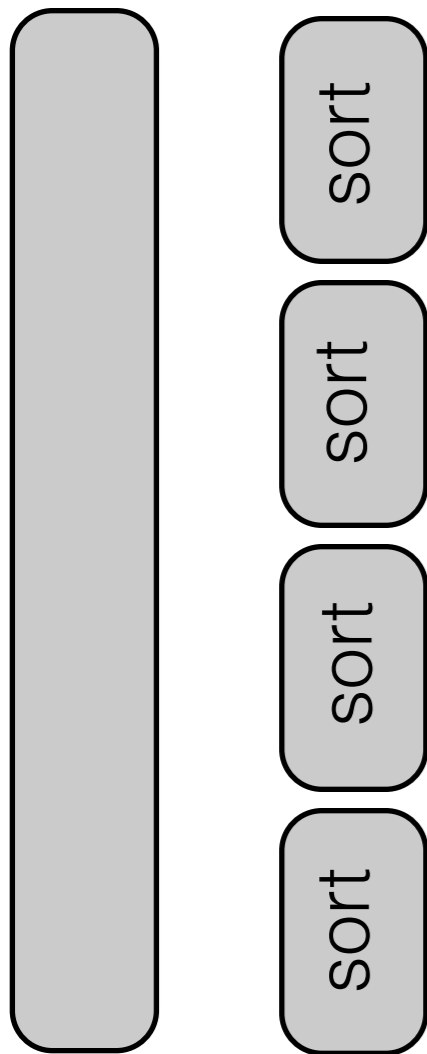
EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`

`select(A,55,70)`

`select(A,150,170)`



adaptive merging

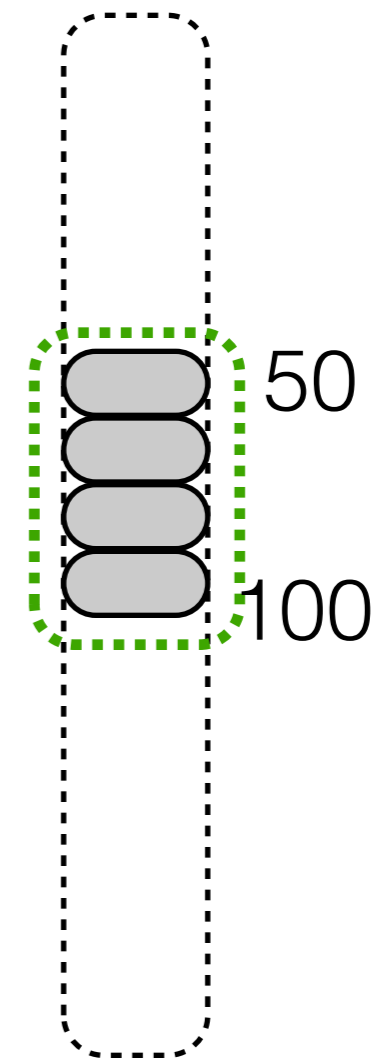
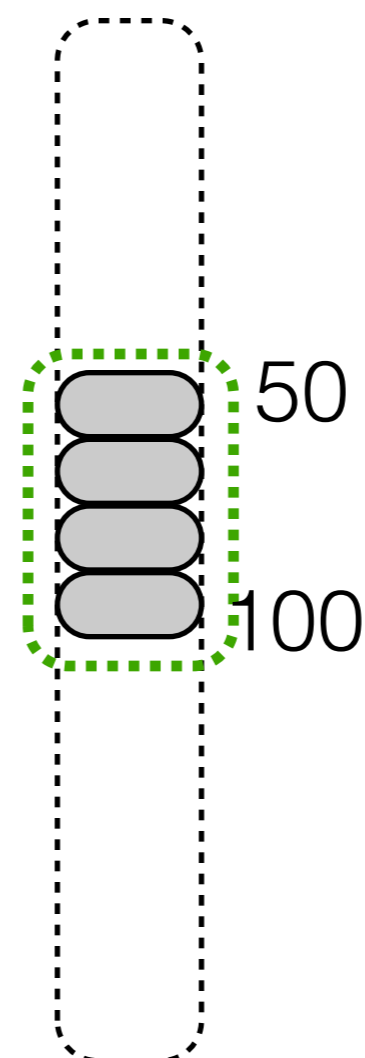
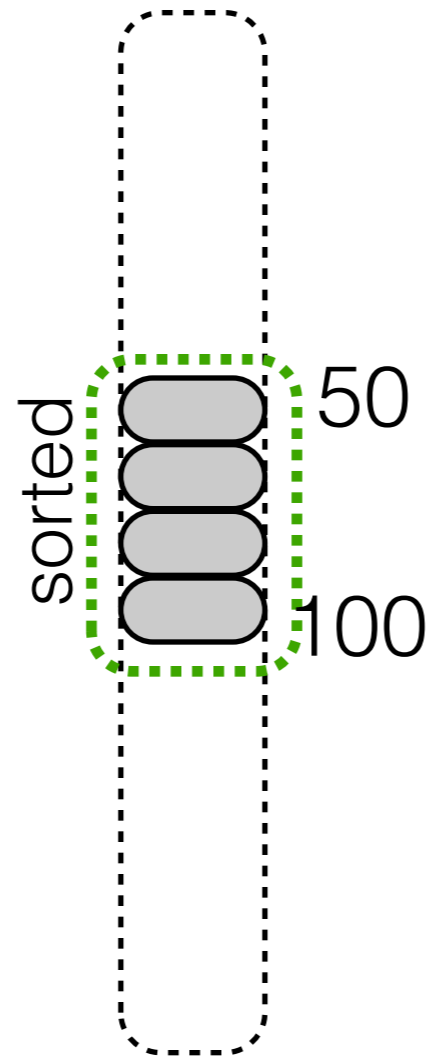
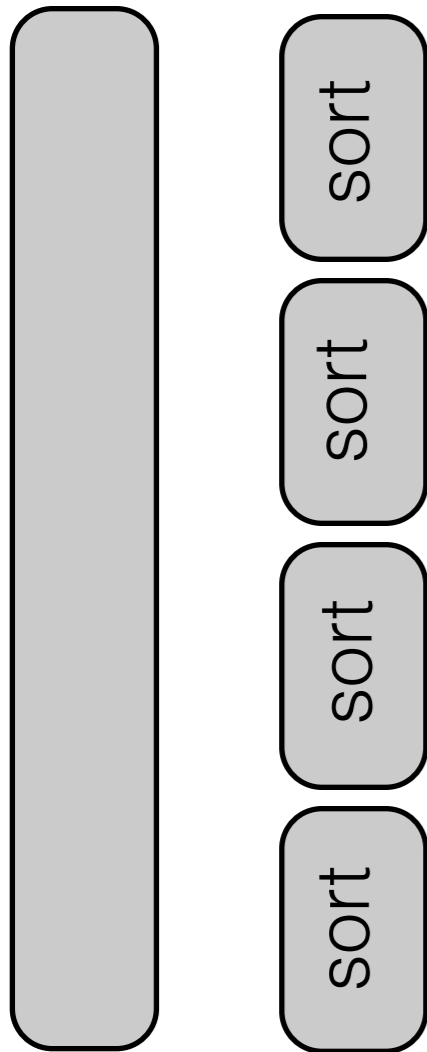
EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

`select(A,50,100)`

`select(A,55,70)`

`select(A,150,170)`



adaptive merging

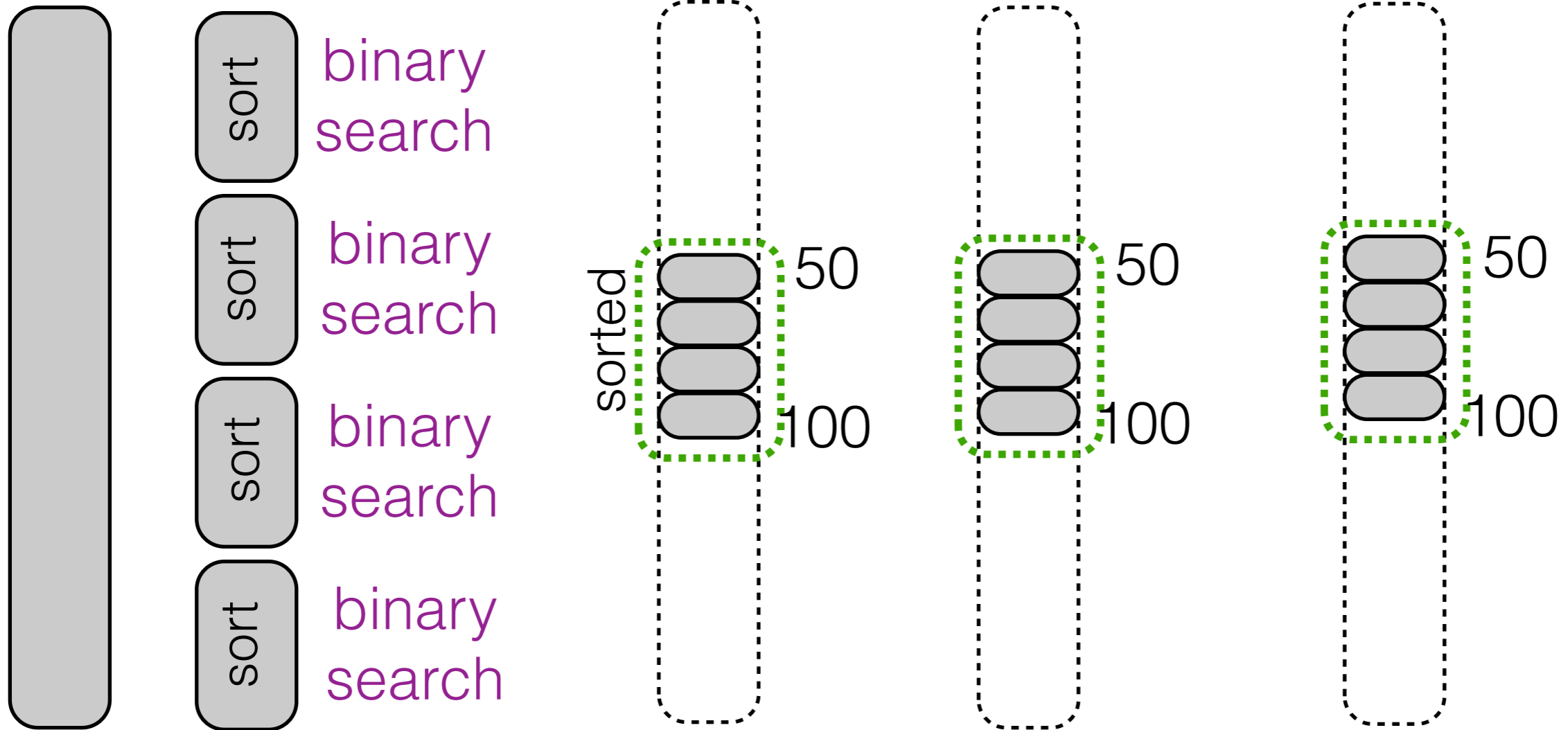
EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

select(A,50,100)

select(A,55,70)

select(A,150,170)



adaptive merging

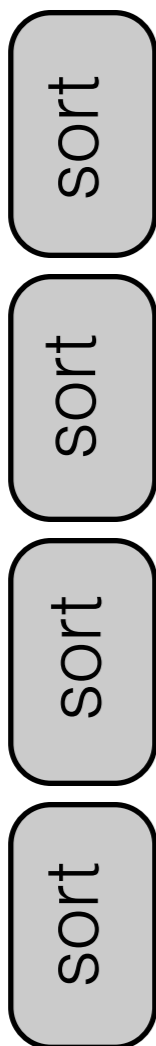
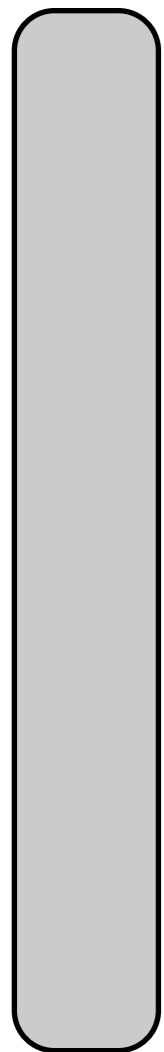
EDBT'10 Goetz Graefe and Harumi Kuno

Incremental sort via external merge sort steps

select(A, 50, 100)

select(A, 55, 70)

select(A, 150, 170)

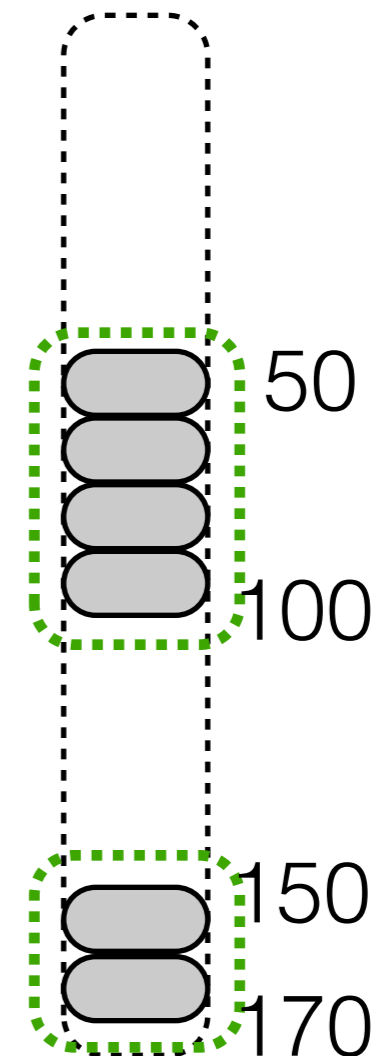
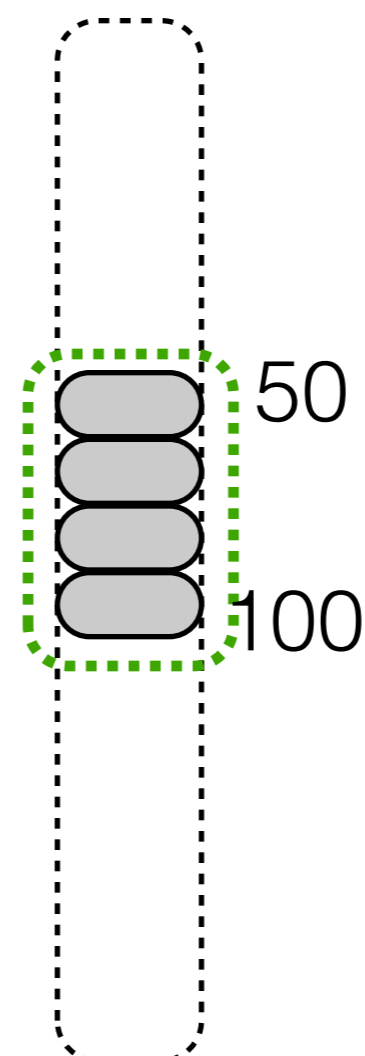
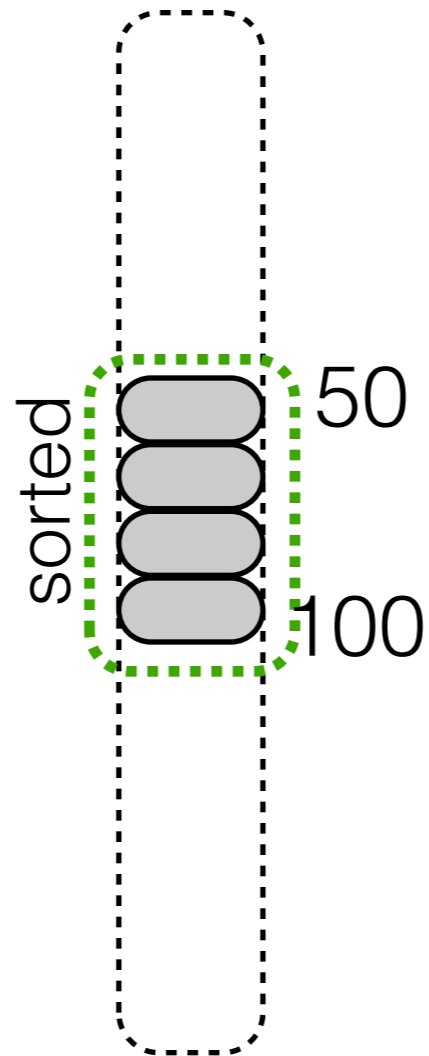


binary search

binary search

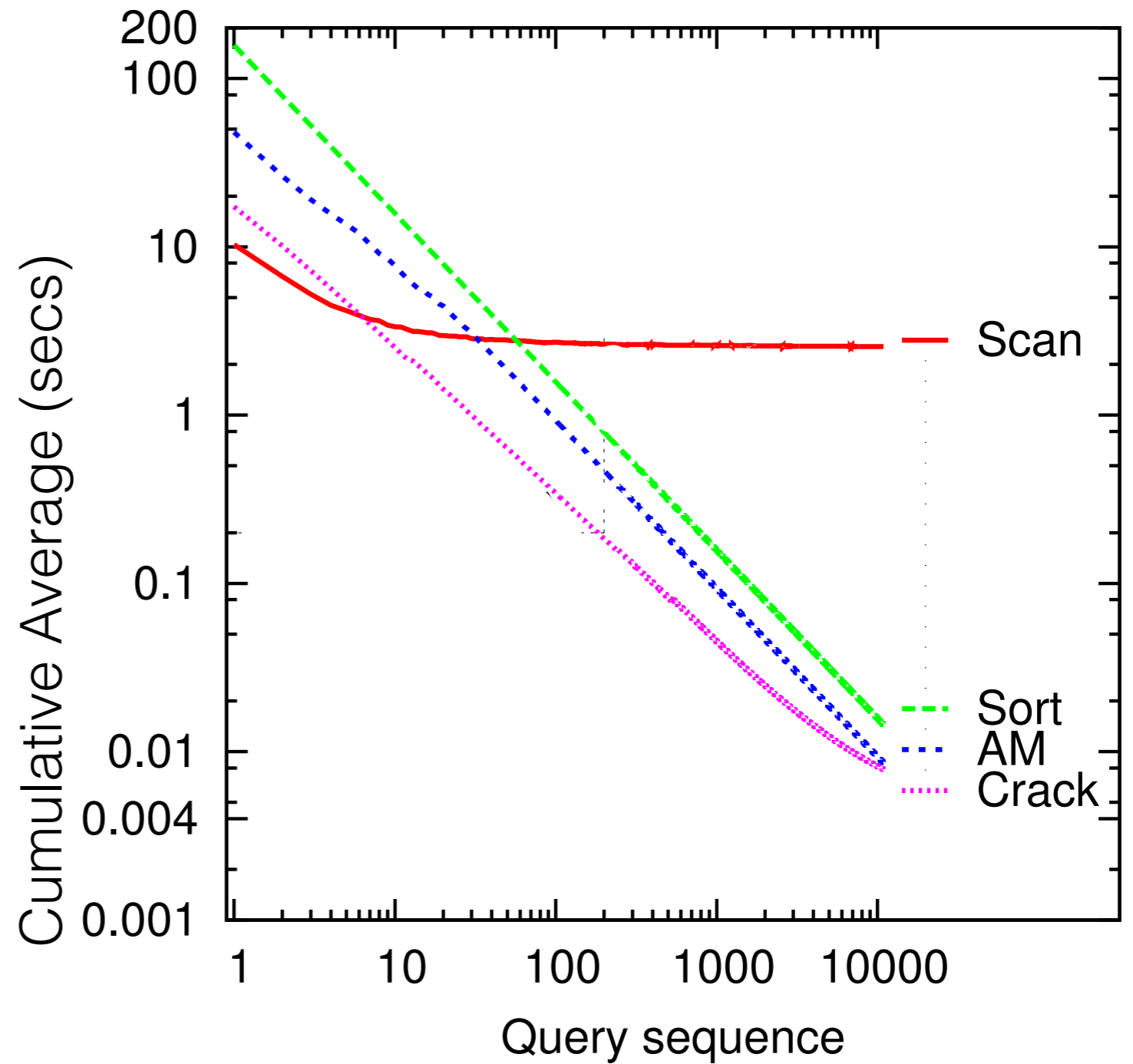
binary search

binary search



set-up

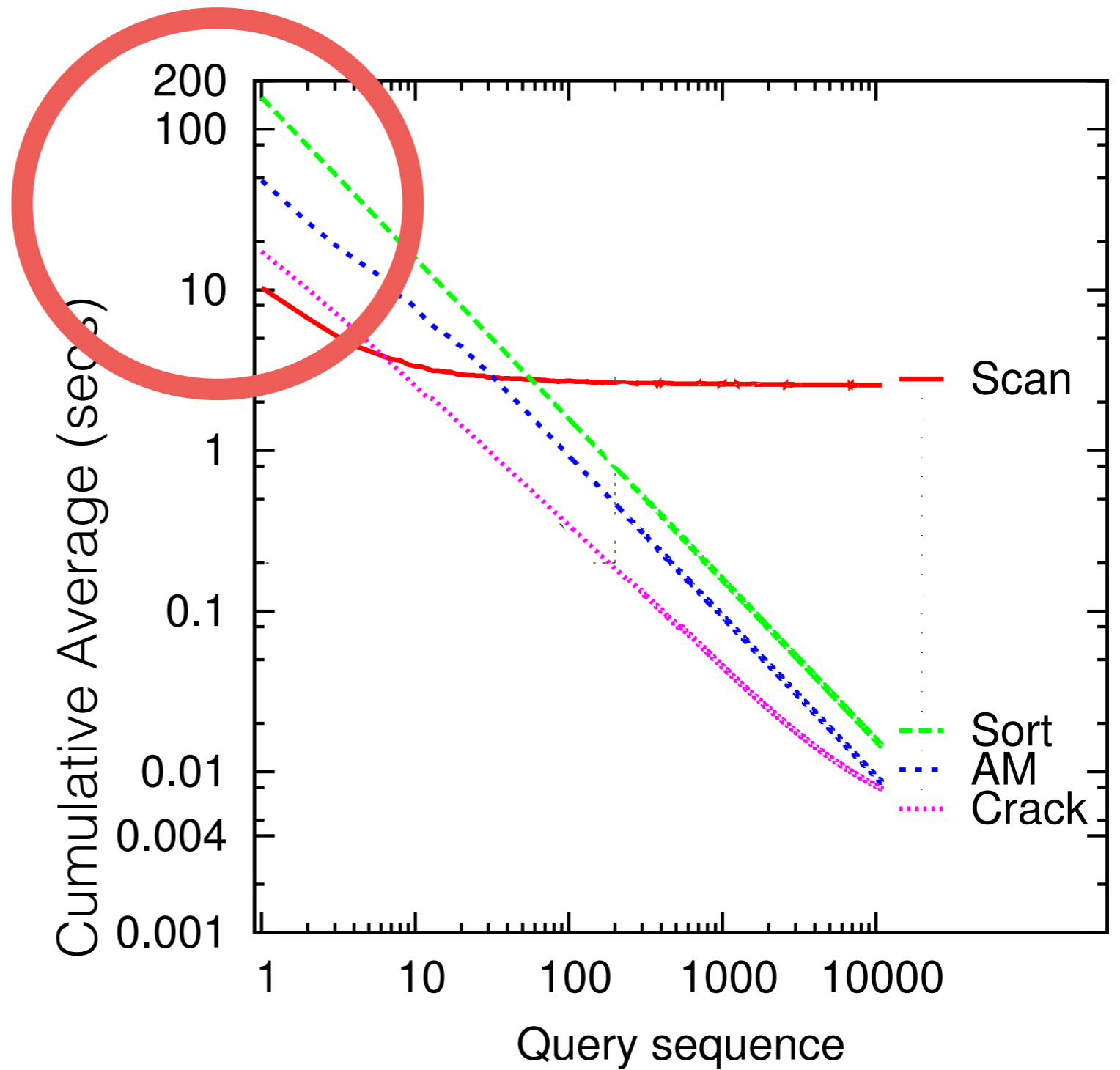
10K random selections
selectivity 10%
random value ranges
in a 30 million integer column



set-up

10K random selections
selectivity 10%
random value ranges
in a 30 million integer column

AM: high init overhead
but fast convergence

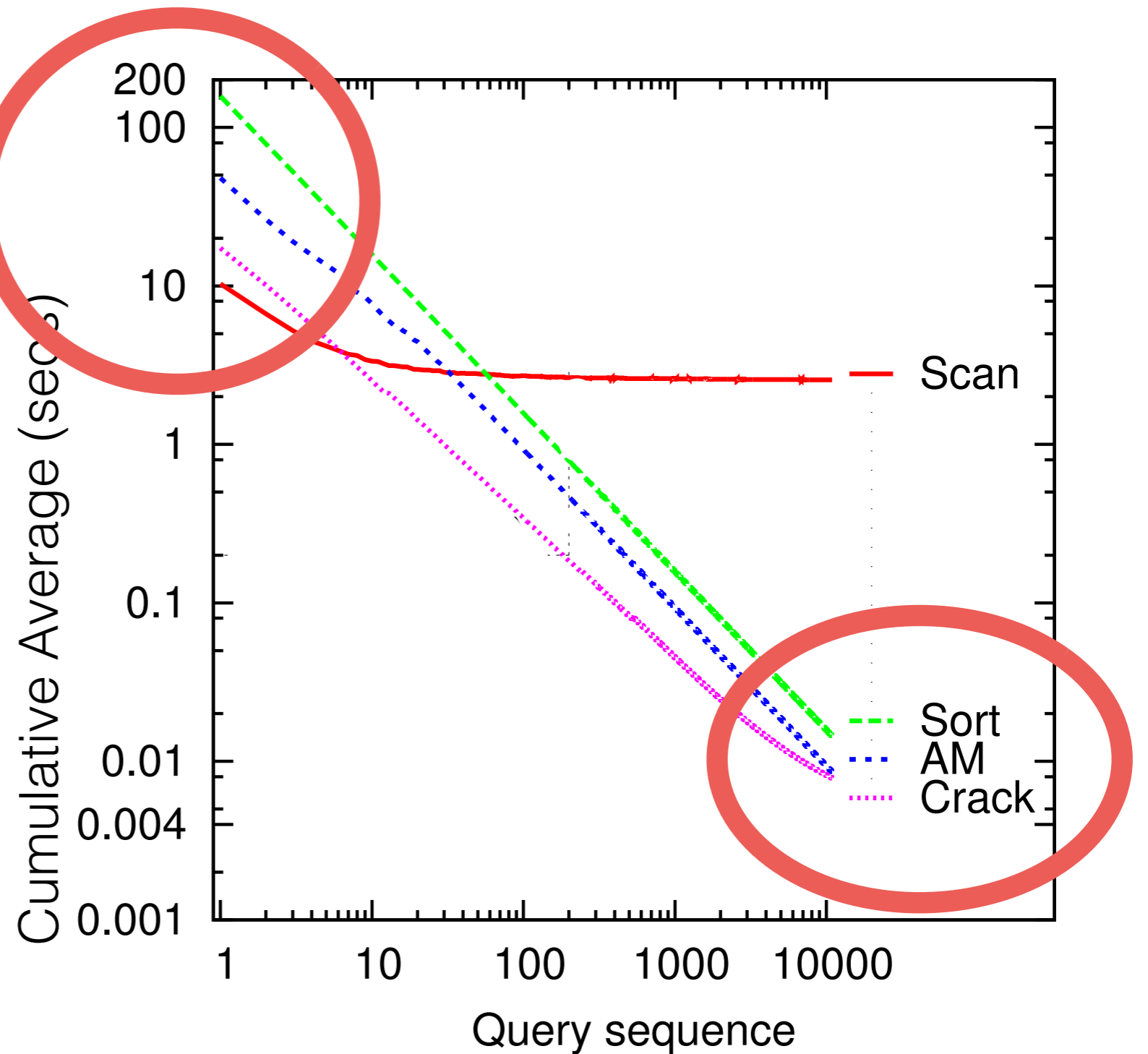


set-up

10K random selections
selectivity 10%
random value ranges
in a 30 million integer column

AM: high init overhead
but fast convergence

Crack: low init overhead
but slow convergence



adaptive merging and cracking are extremes

what is there in between?

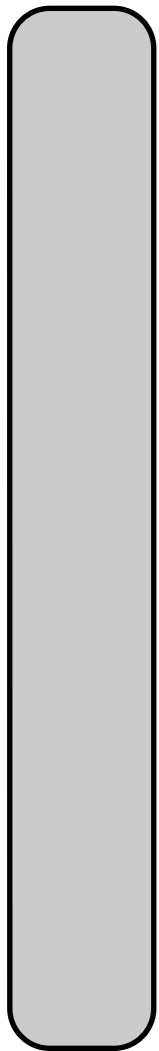
crack-crack

vary initialization and incremental steps taken

crack-crack

vary initialization and incremental steps taken

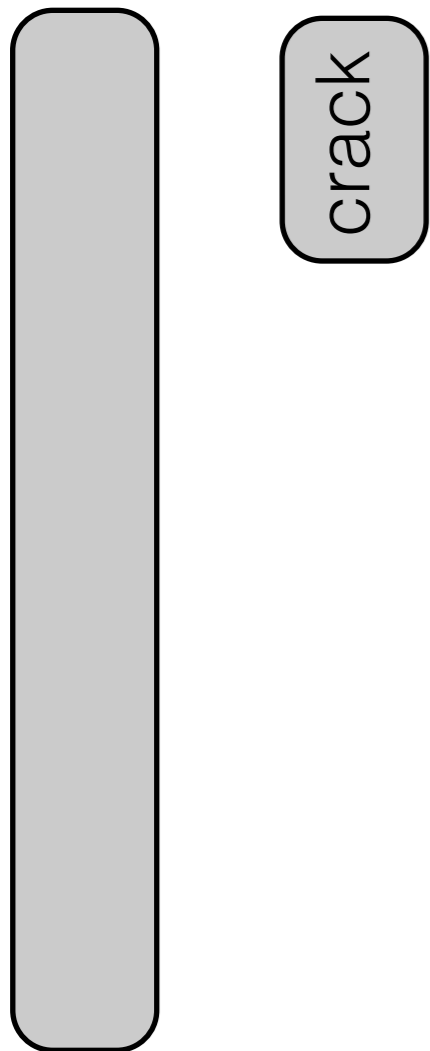
`select(A,50,100)`



crack-crack

vary initialization and incremental steps taken

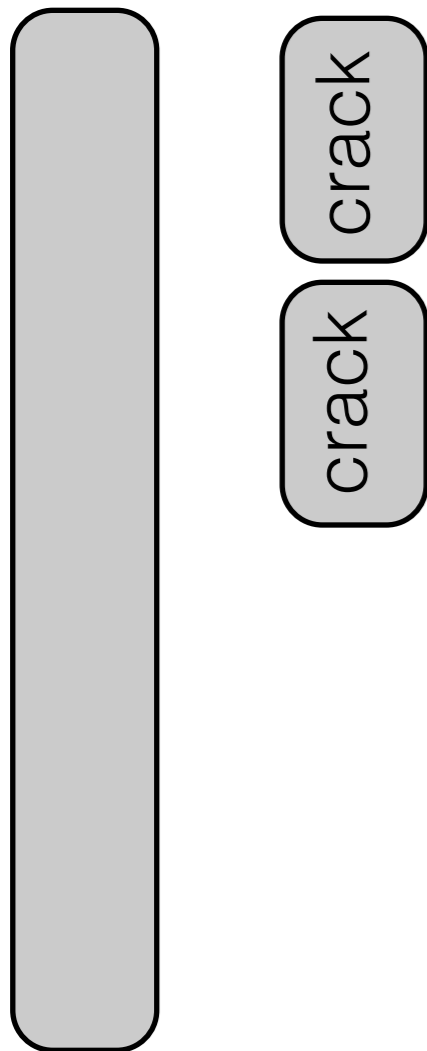
`select(A,50,100)`



crack-crack

vary initialization and incremental steps taken

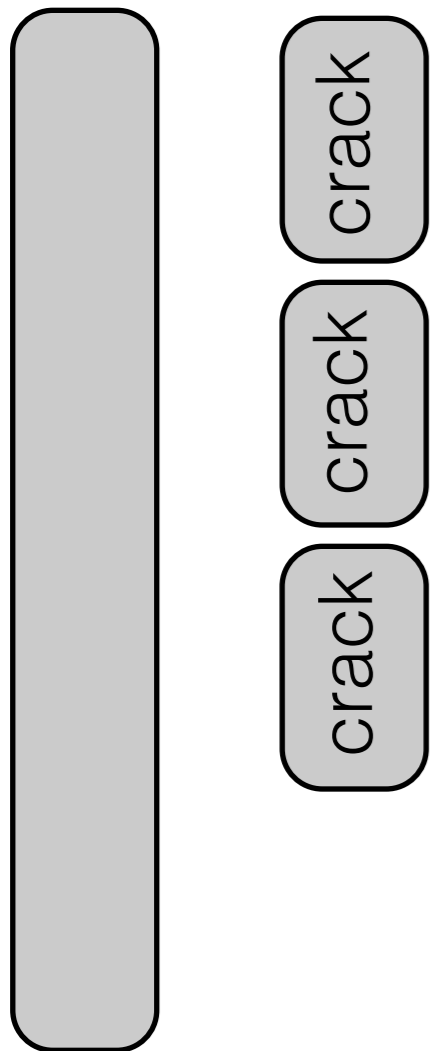
`select(A, 50, 100)`



crack-crack

vary initialization and incremental steps taken

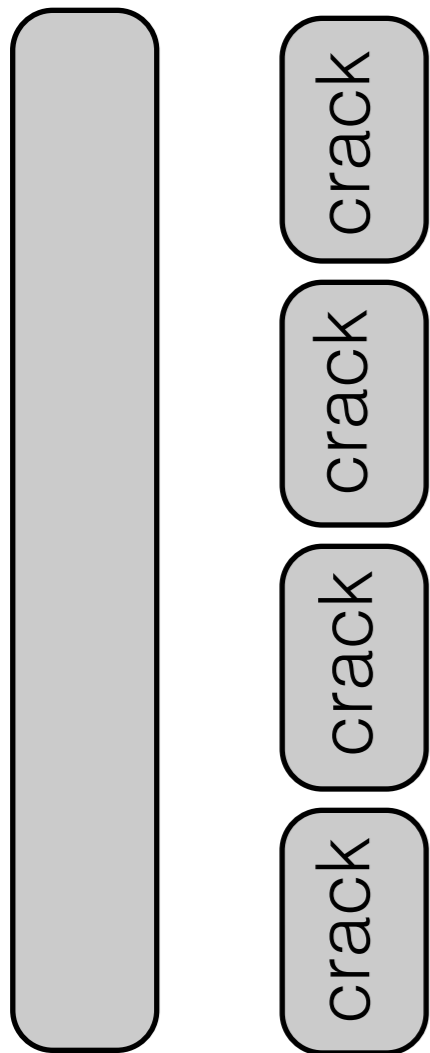
`select(A,50,100)`



crack-crack

vary initialization and incremental steps taken

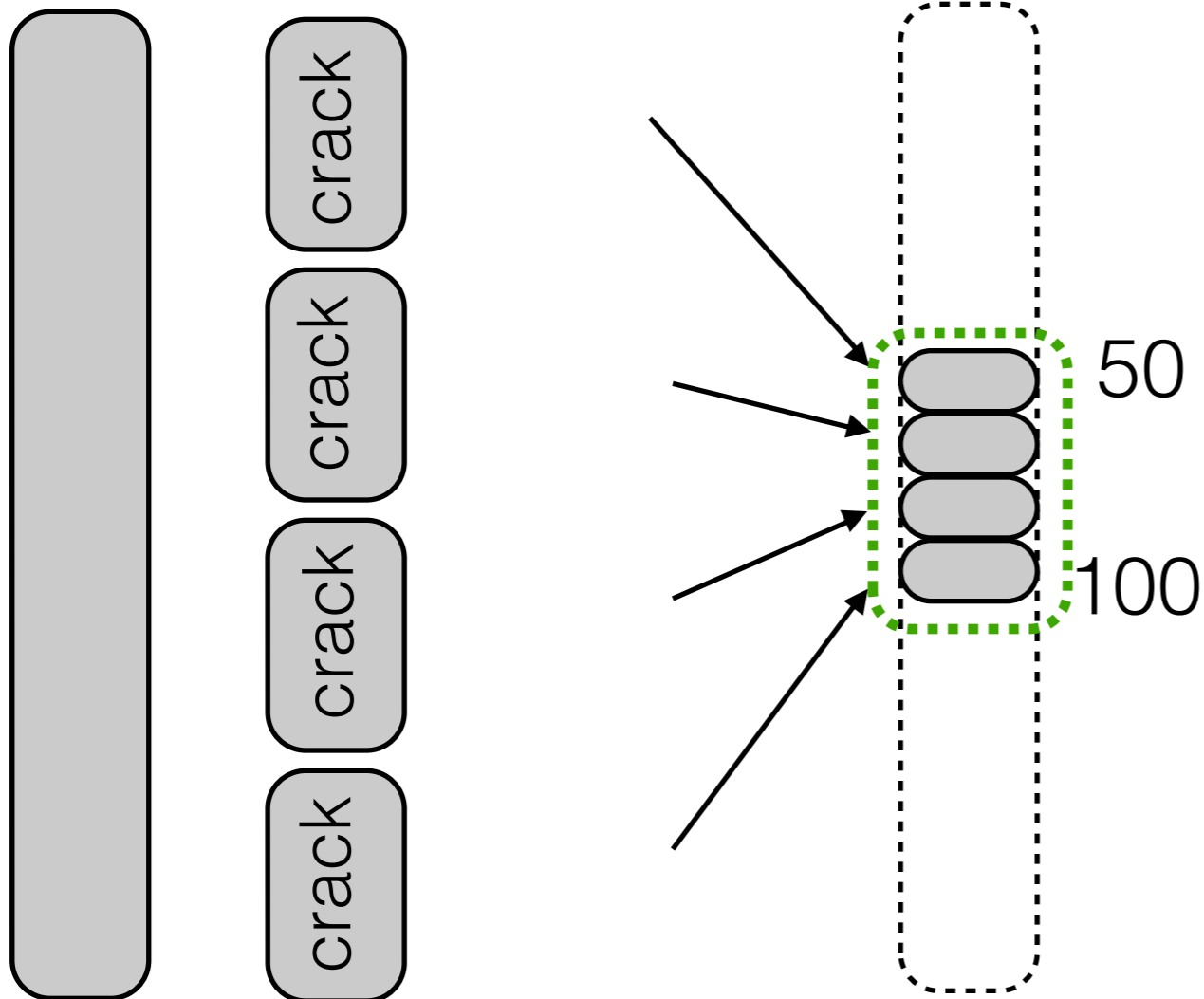
`select(A, 50, 100)`



crack-crack

vary initialization and incremental steps taken

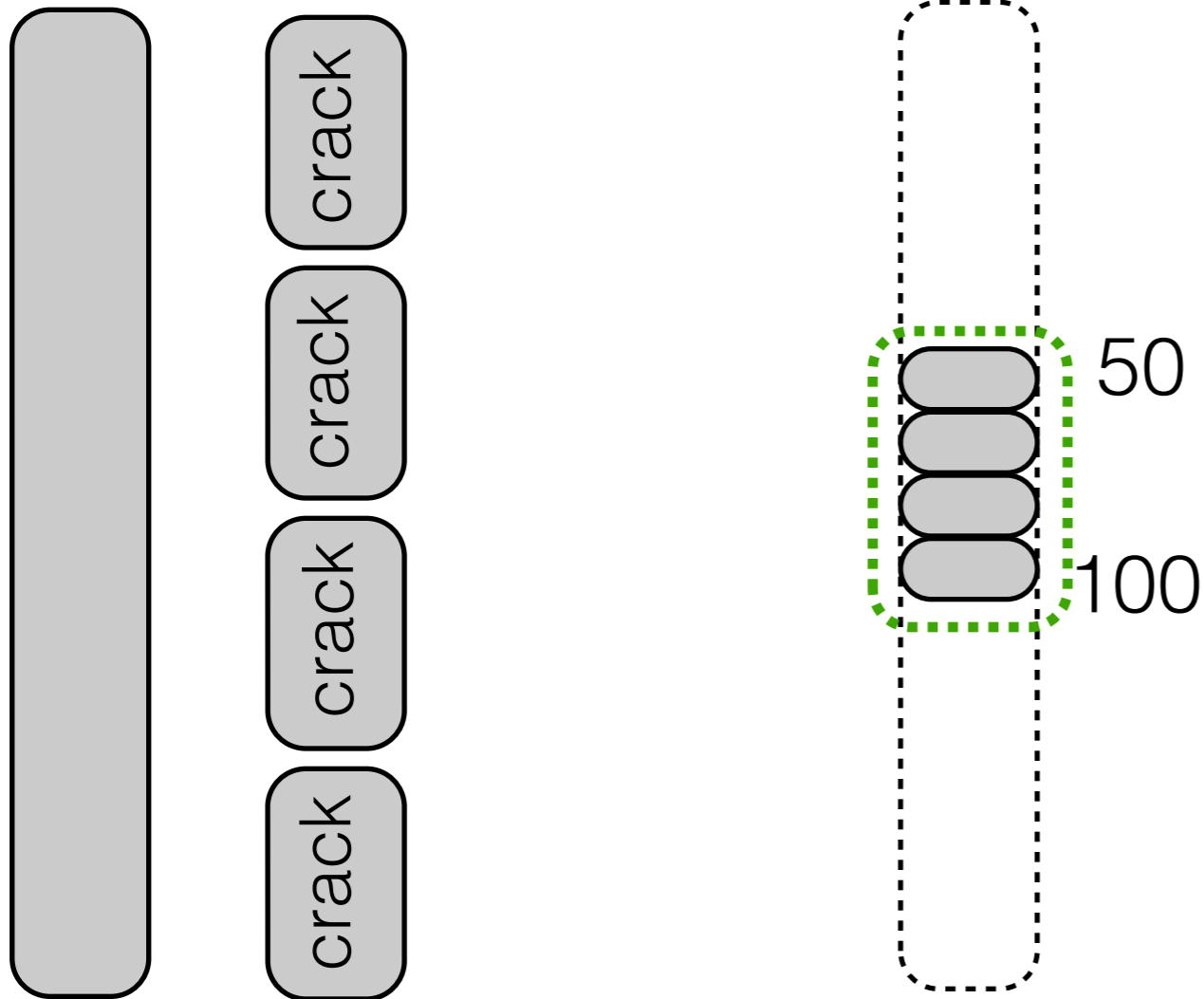
select(A, 50, 100)



crack-crack

vary initialization and incremental steps taken

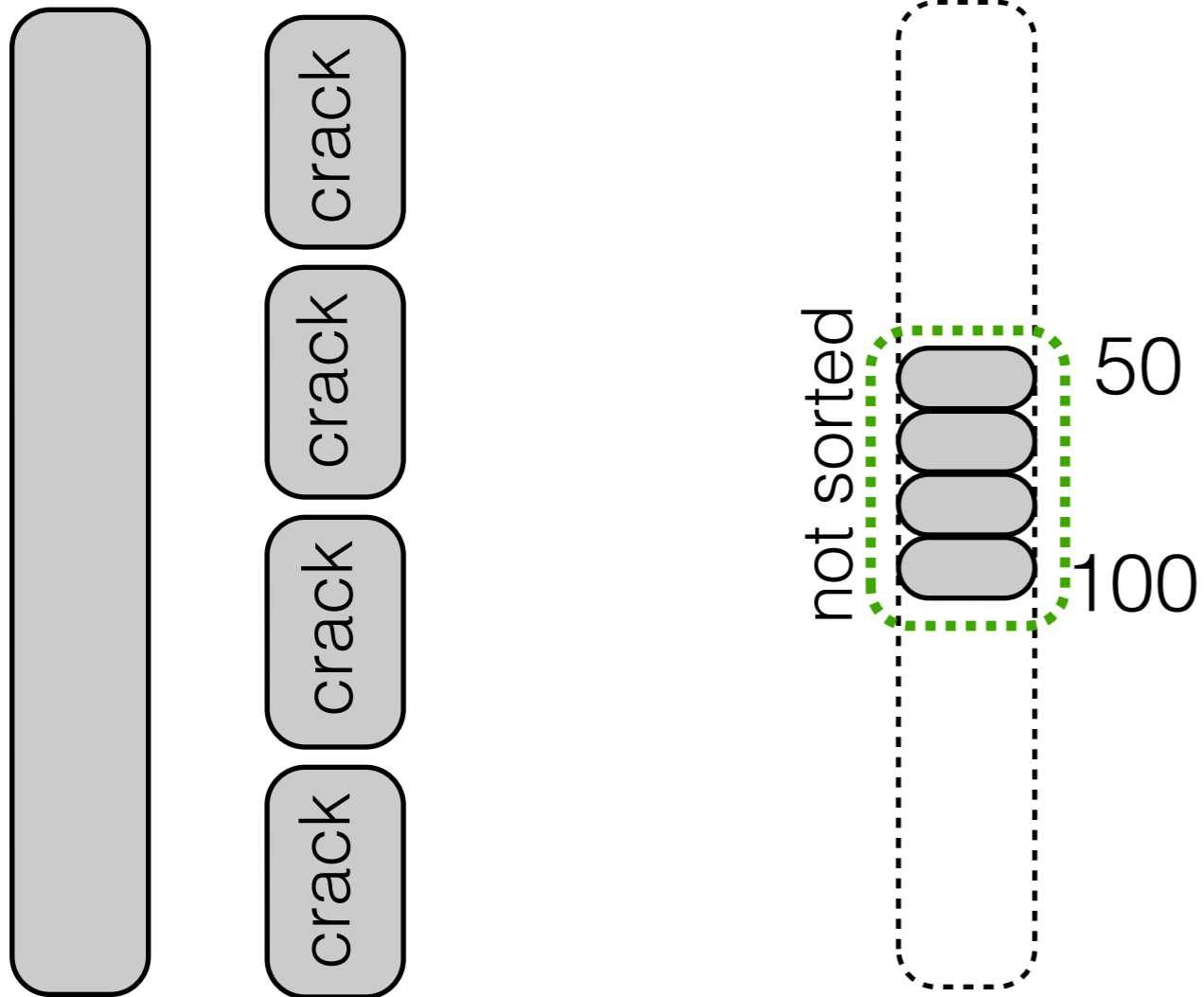
select(A, 50, 100)



crack-crack

vary initialization and incremental steps taken

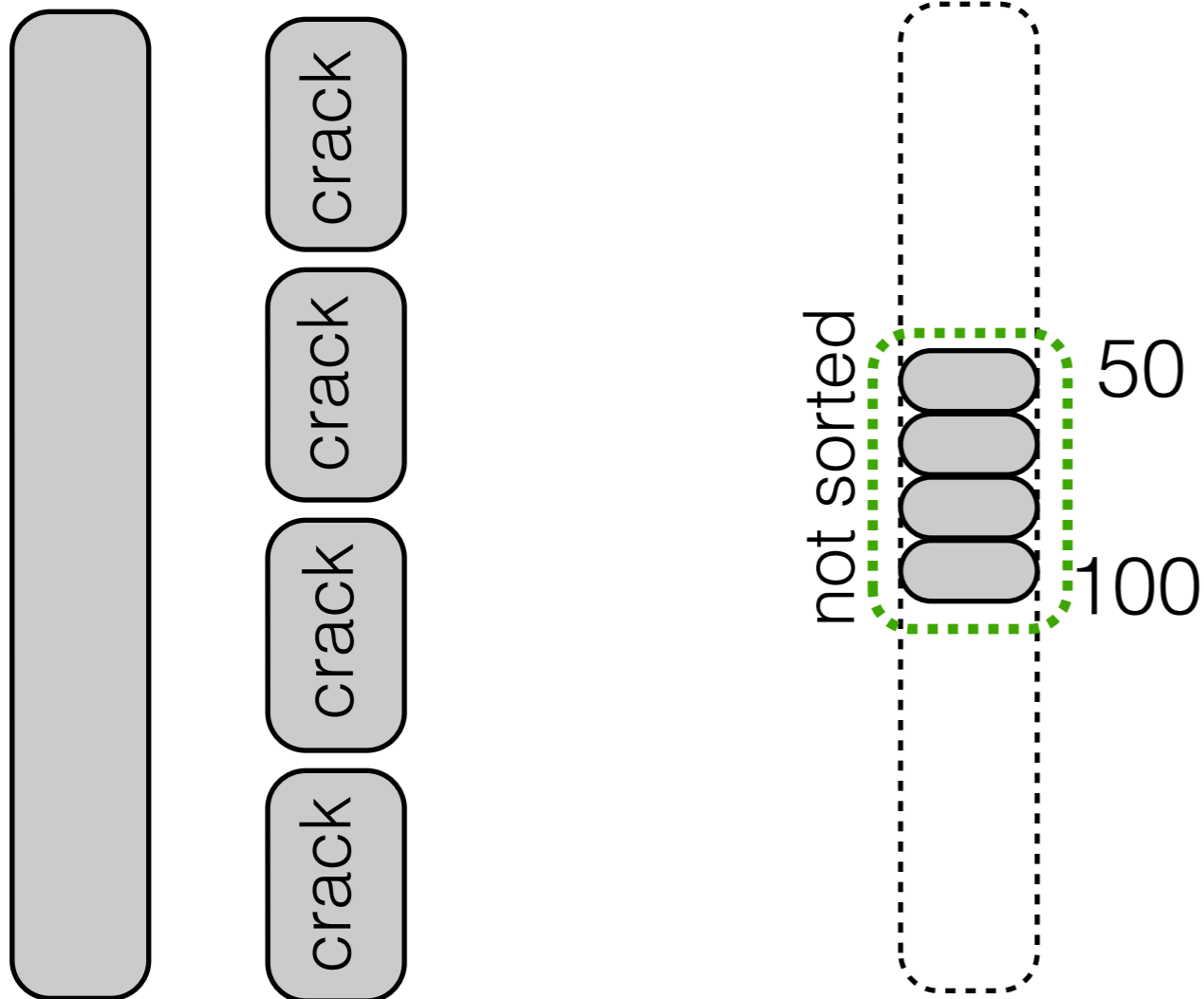
select(A, 50, 100)



crack-crack

vary initialization and incremental steps taken

`select(A,50,100)` `select(A,55,70)`

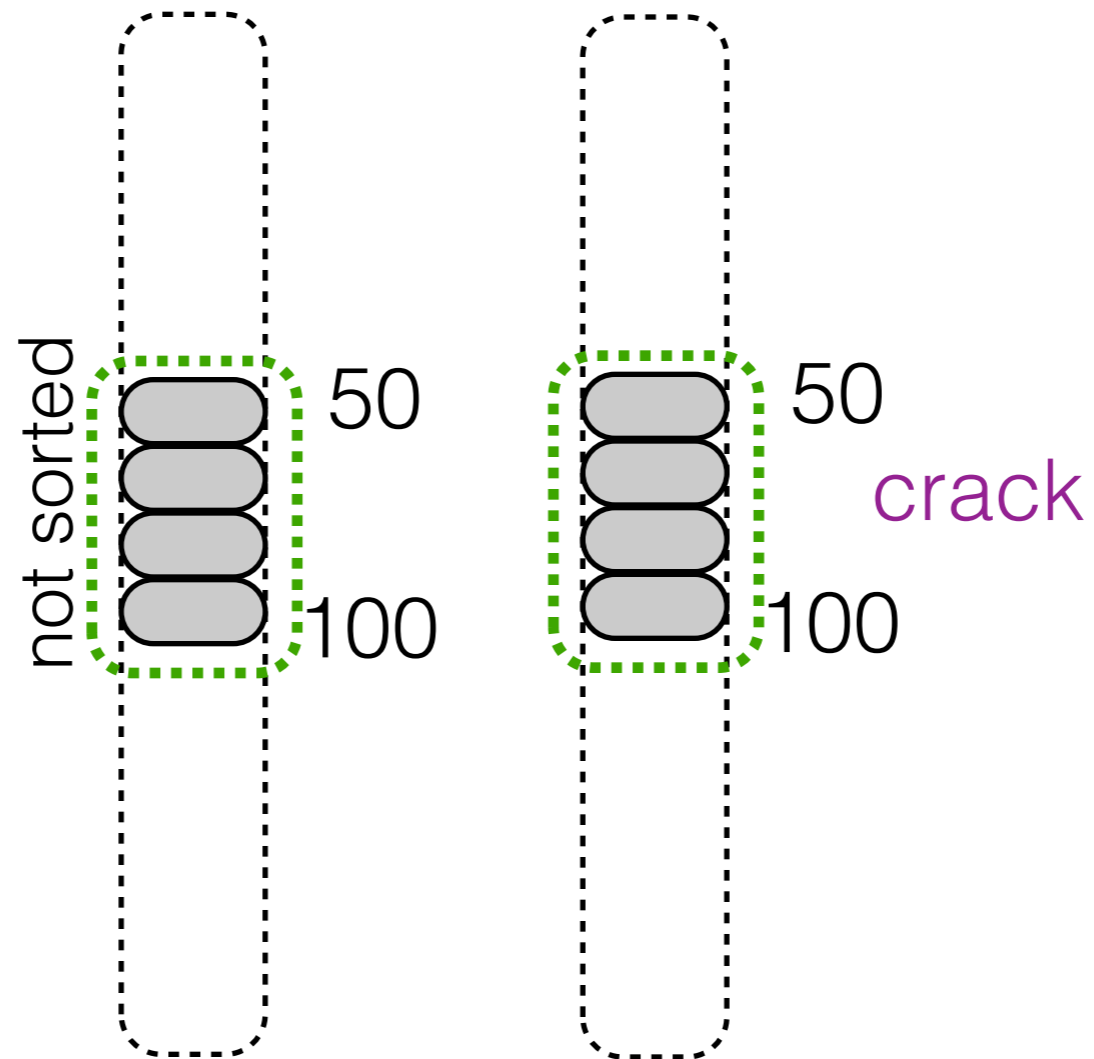
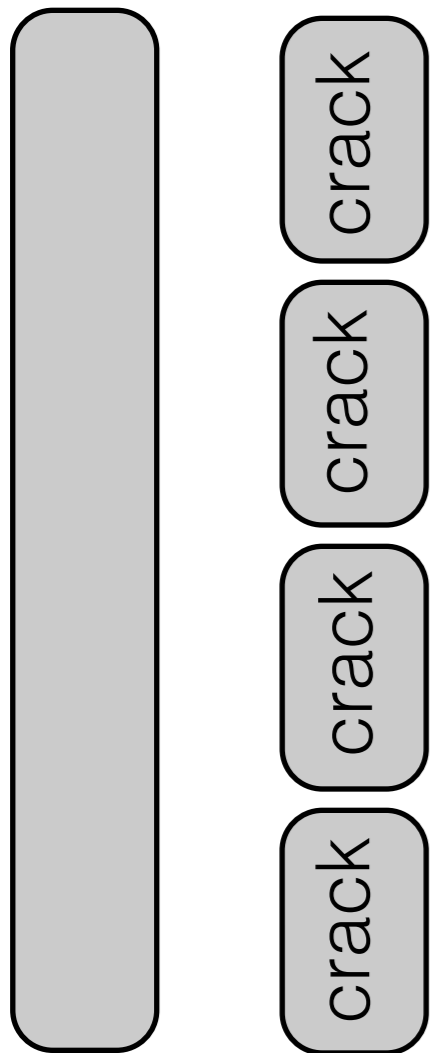


crack-crack

vary initialization and incremental steps taken

`select(A,50,100)`

`select(A,55,70)`



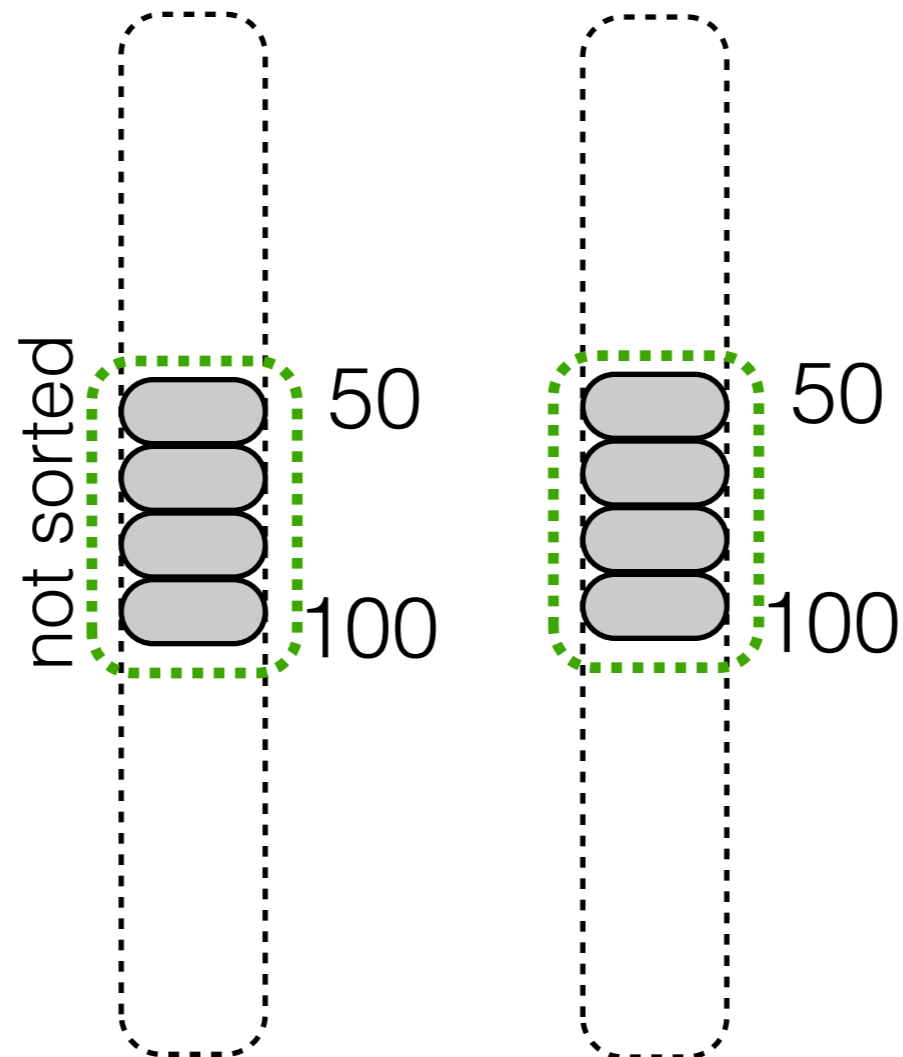
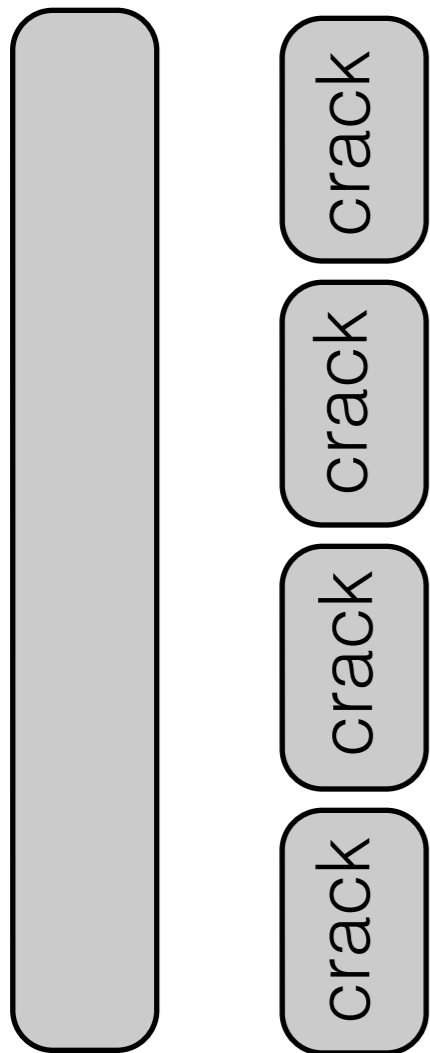
crack-crack

vary initialization and incremental steps taken

select(A, 50, 100)

select(A, 55, 70)

select(A, 150, 170)



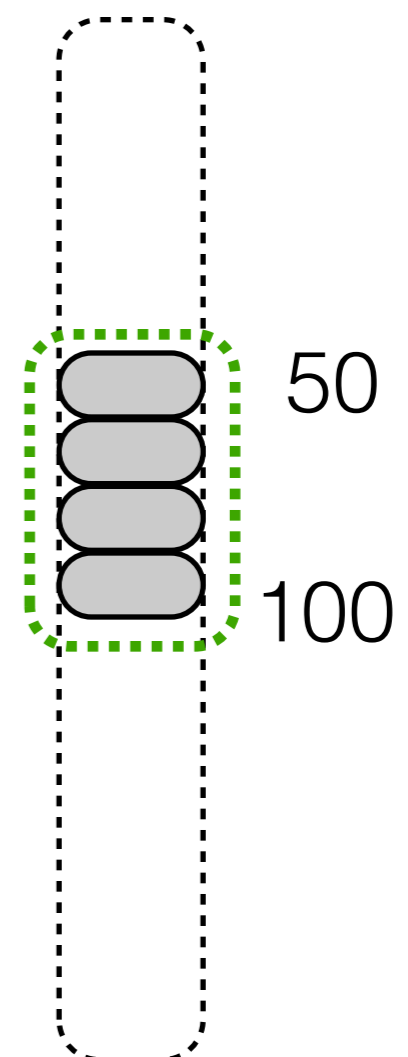
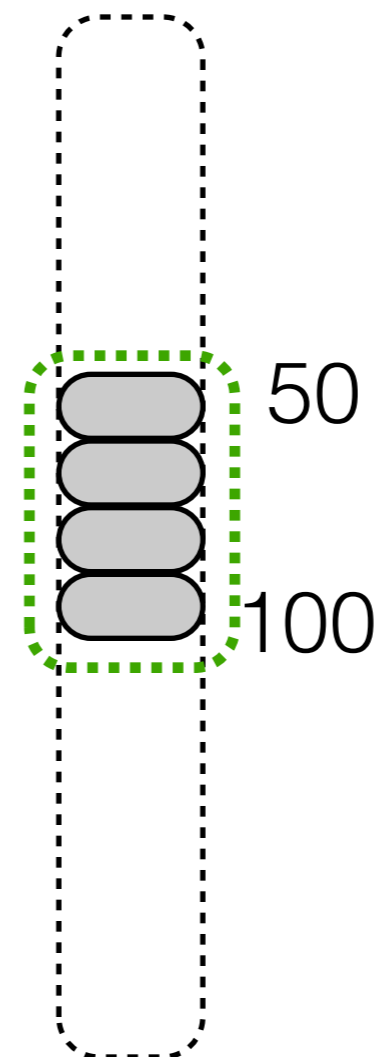
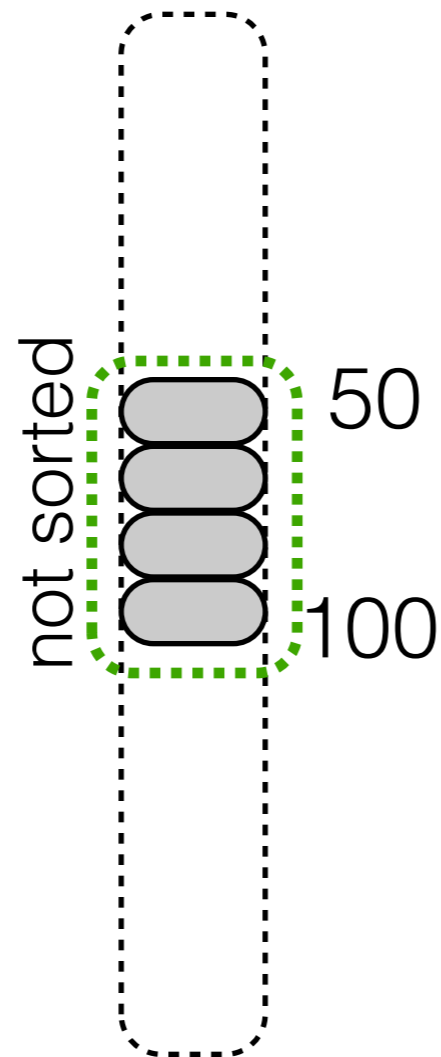
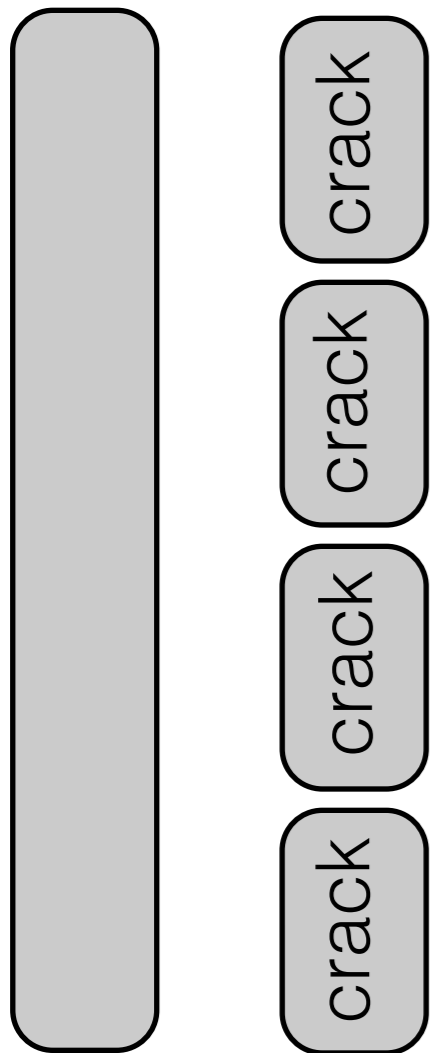
crack-crack

vary initialization and incremental steps taken

`select(A, 50, 100)`

`select(A, 55, 70)`

`select(A, 150, 170)`



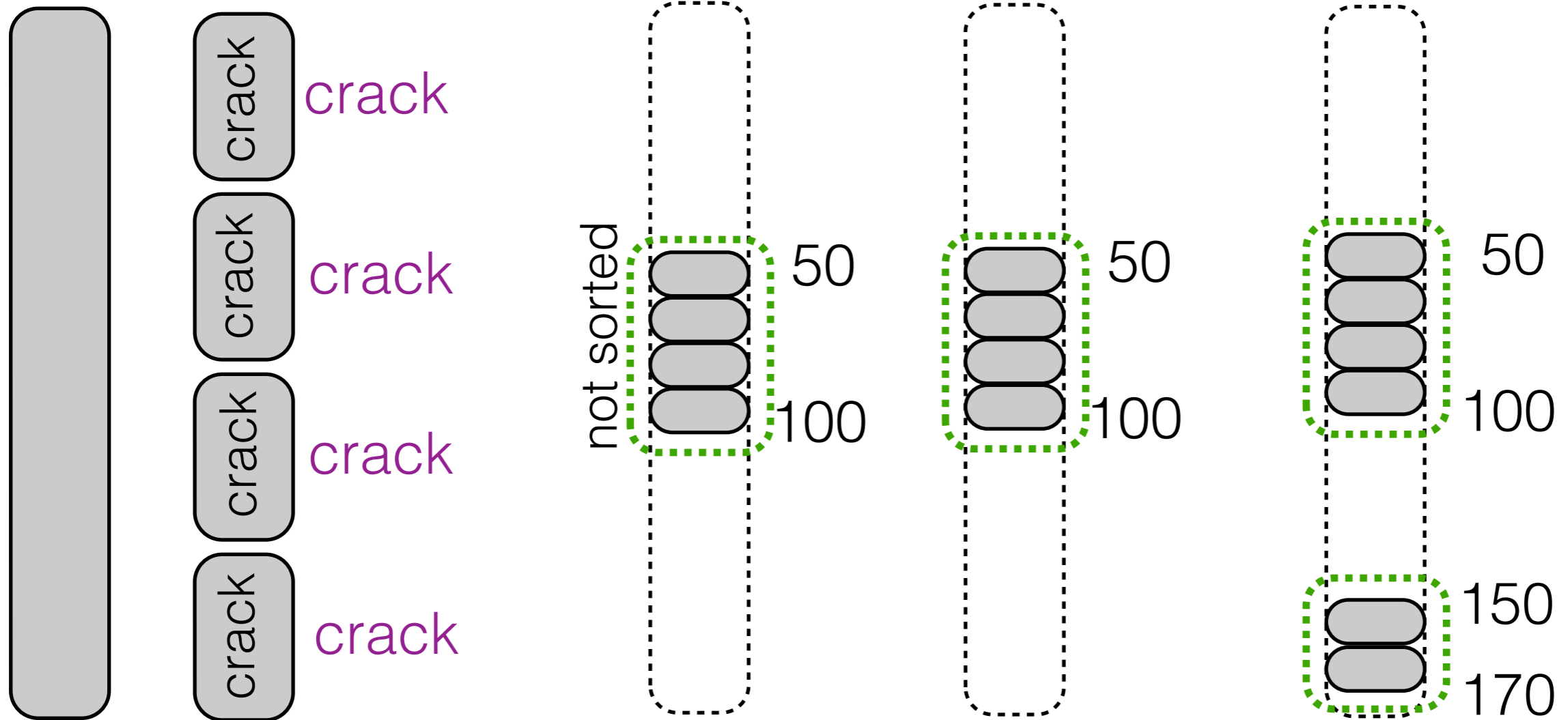
crack-crack

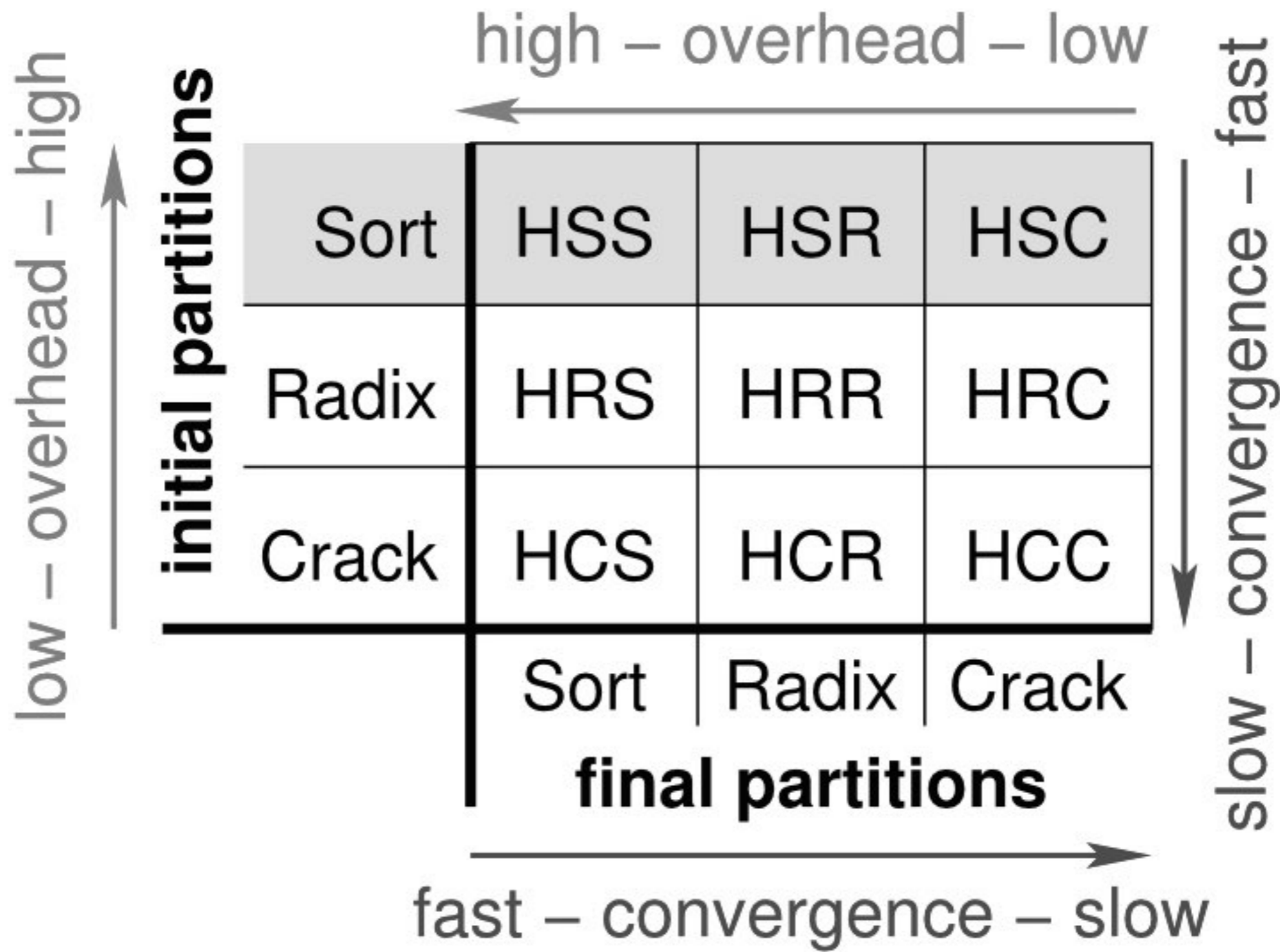
vary initialization and incremental steps taken

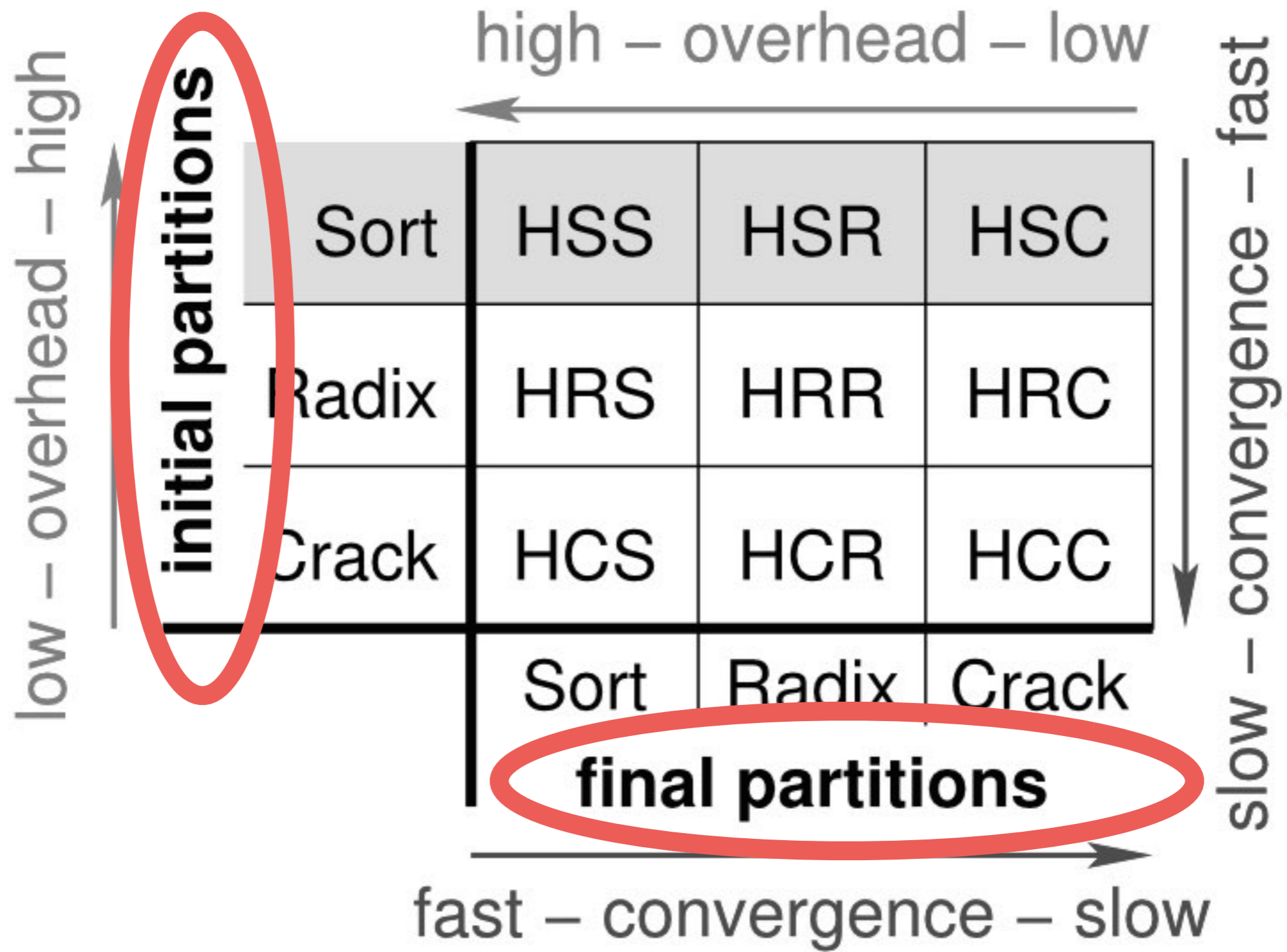
`select(A, 50, 100)`

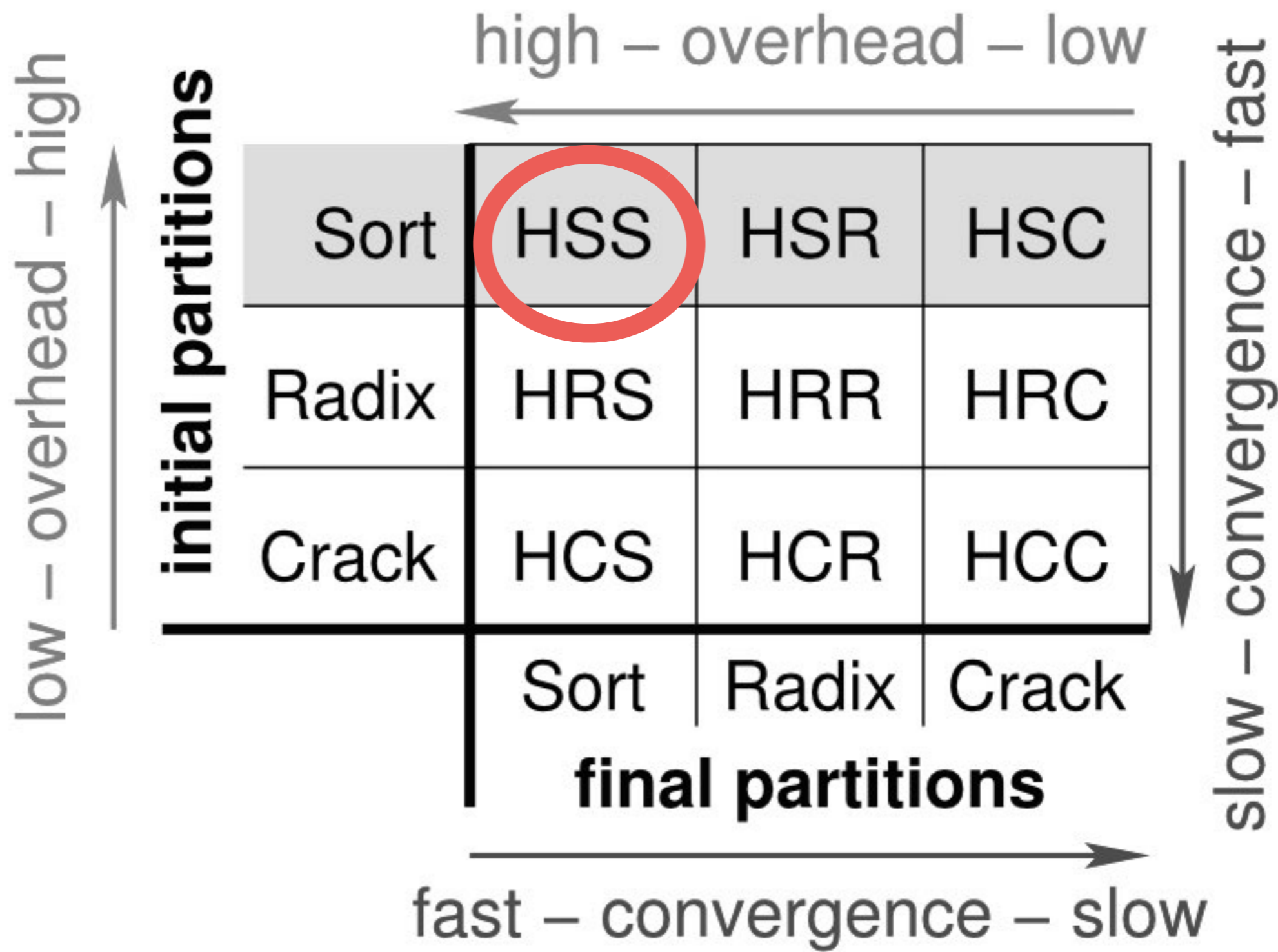
`select(A, 55, 70)`

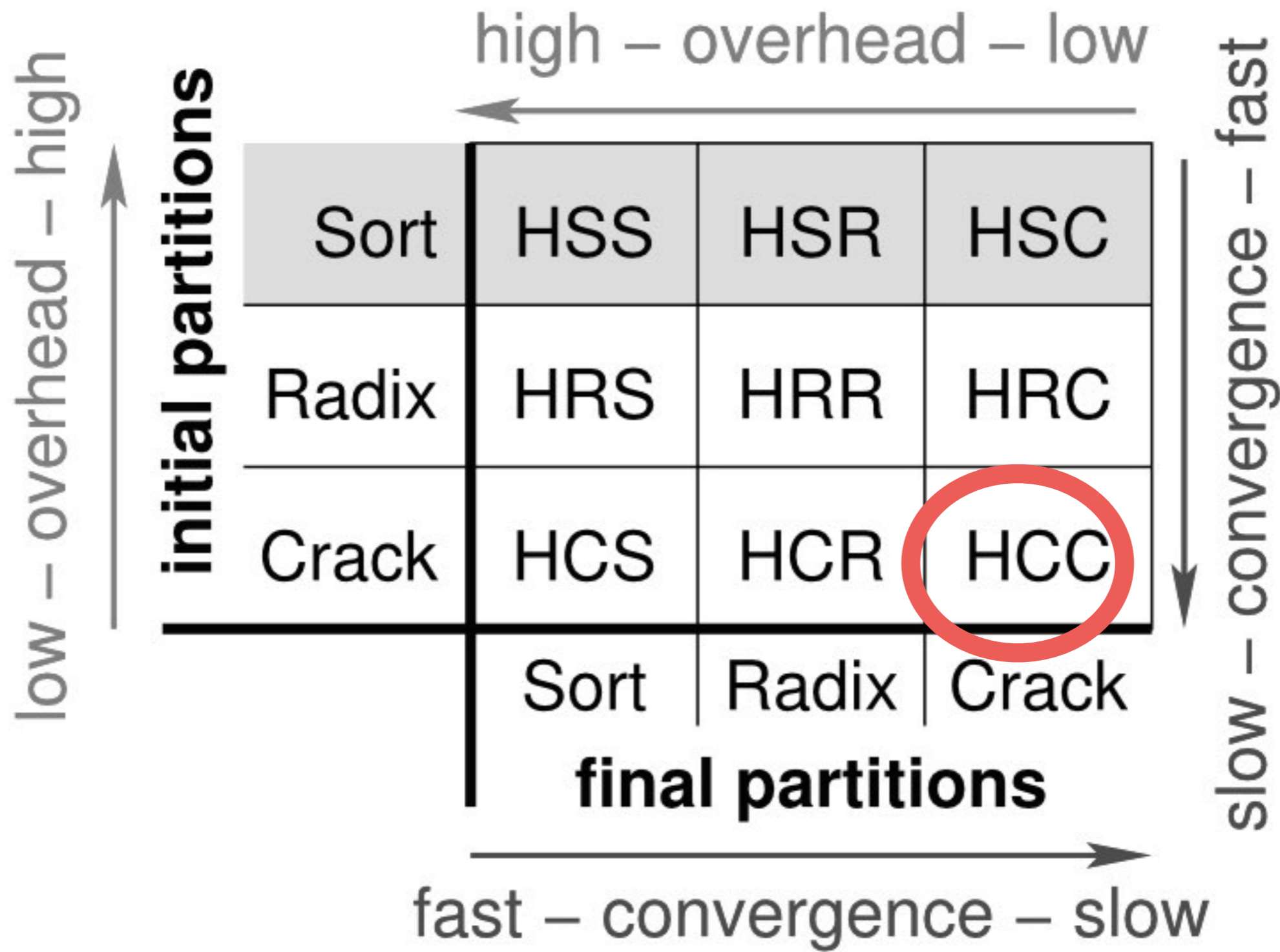
`select(A, 150, 170)`

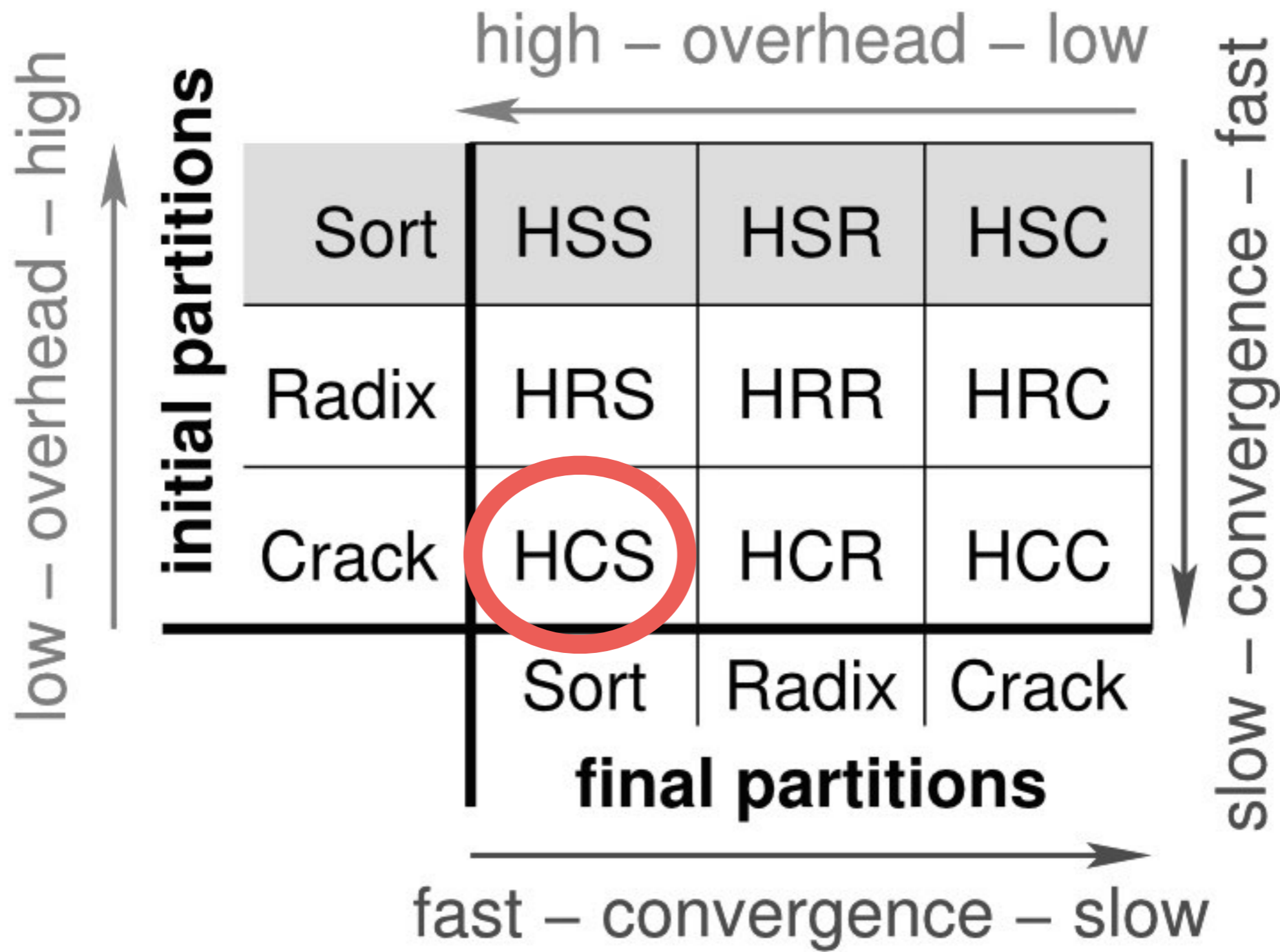


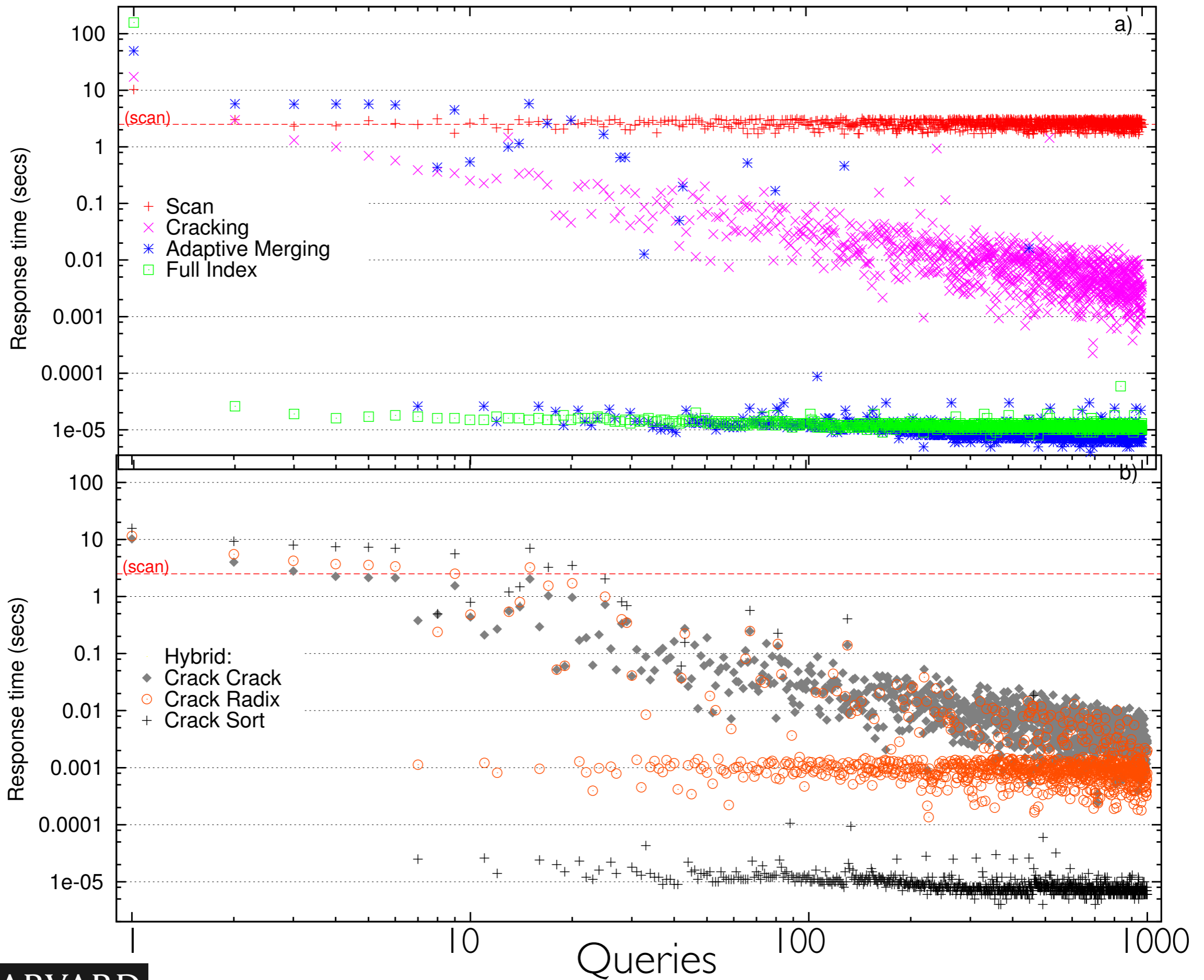


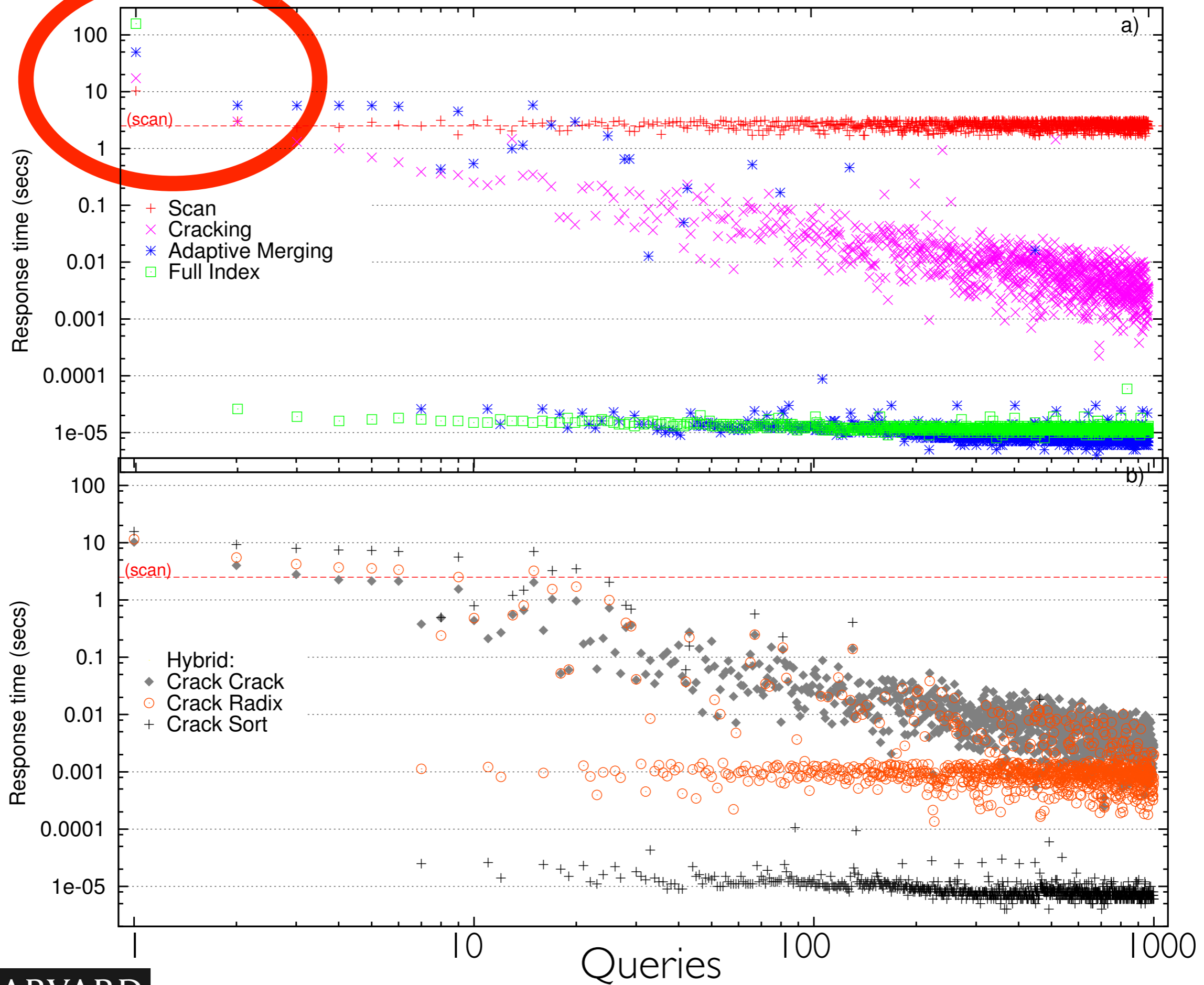


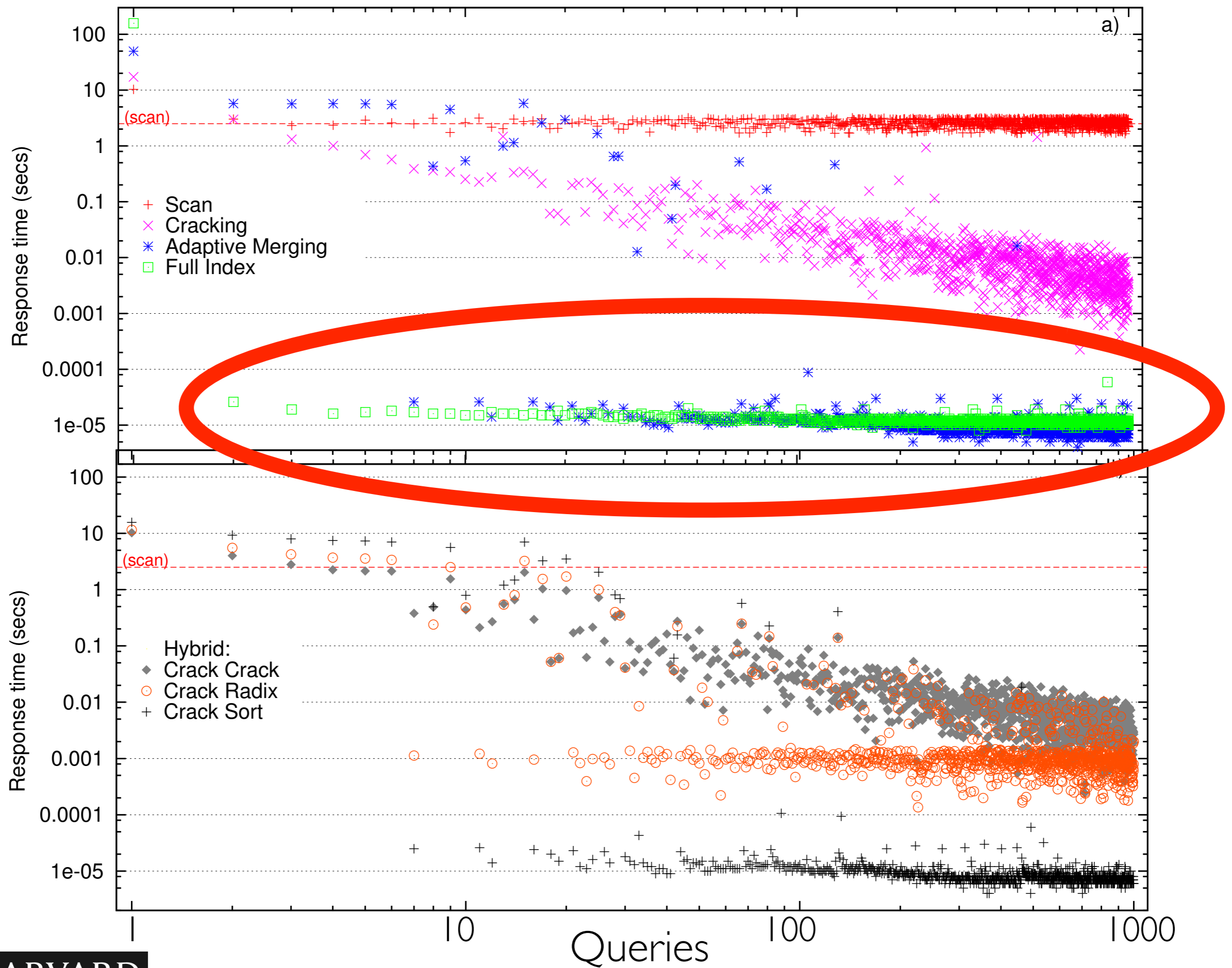


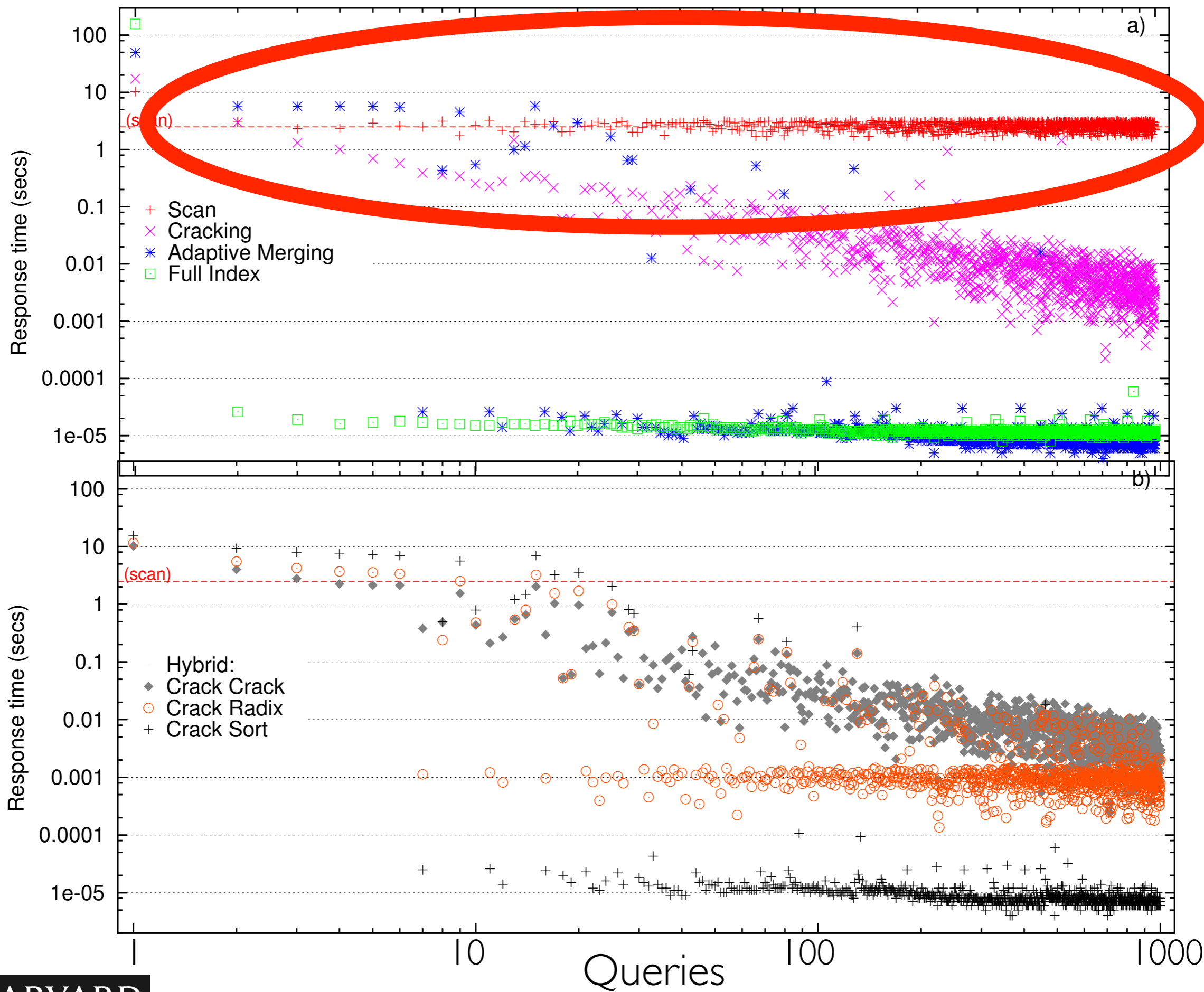


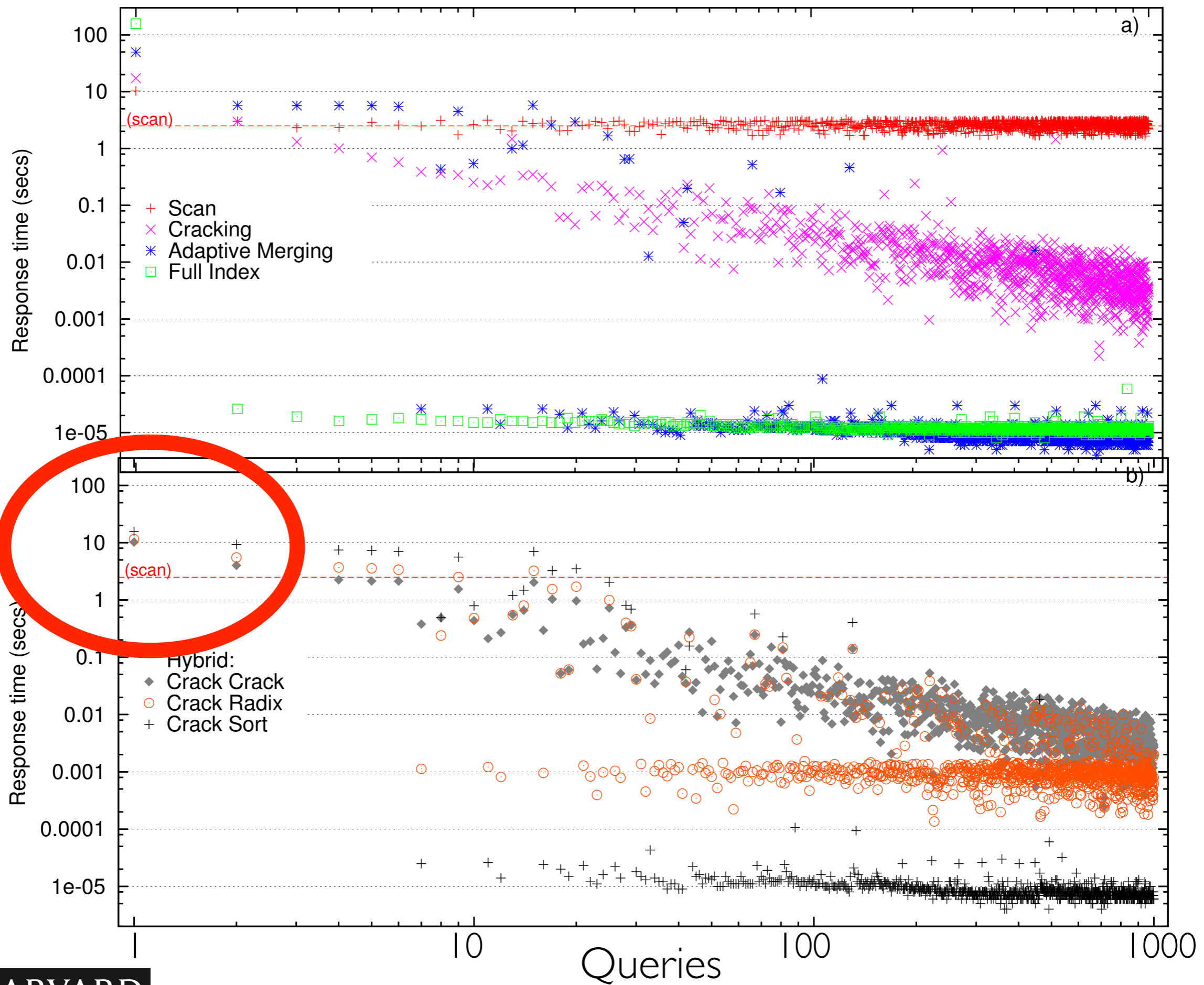


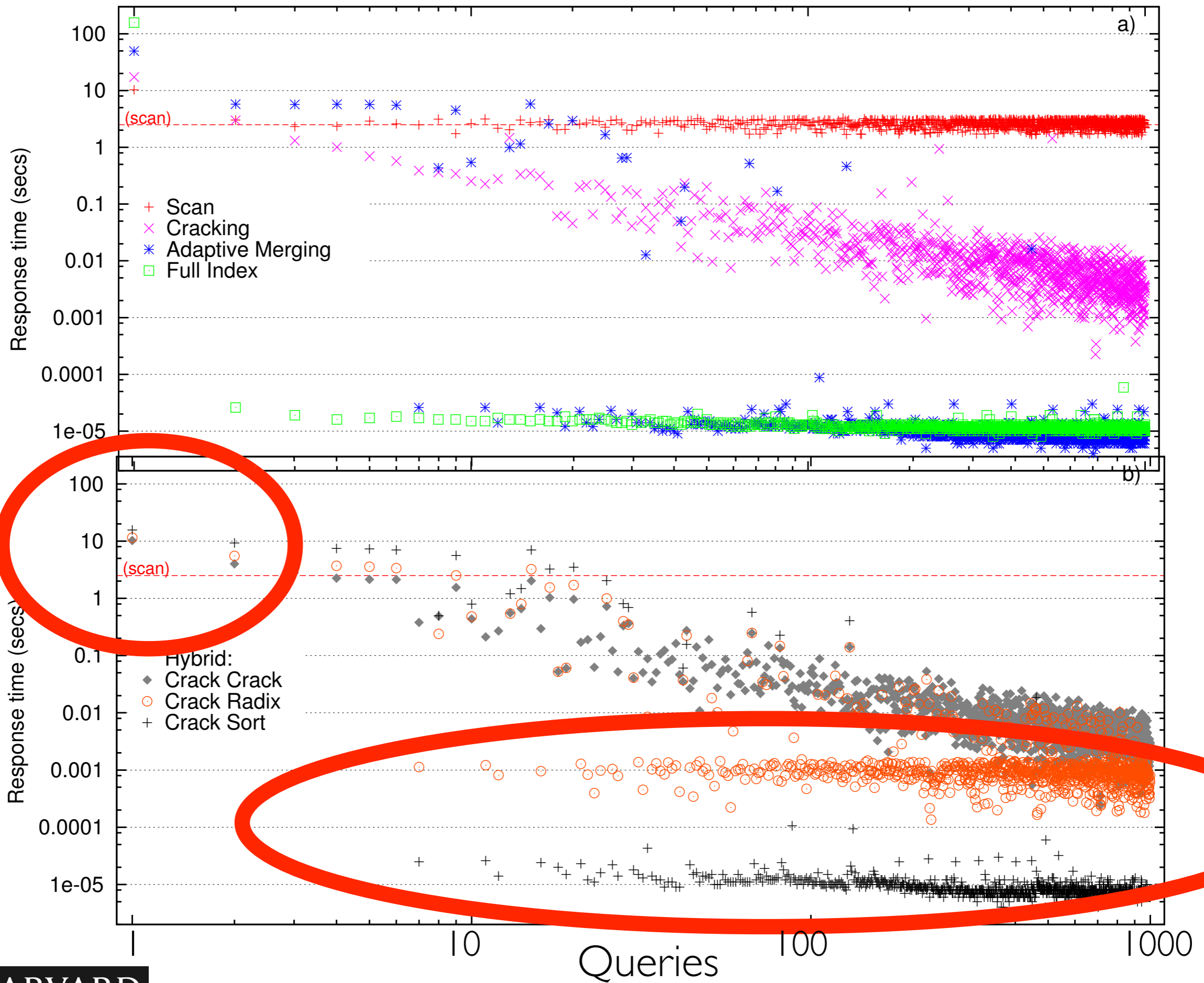




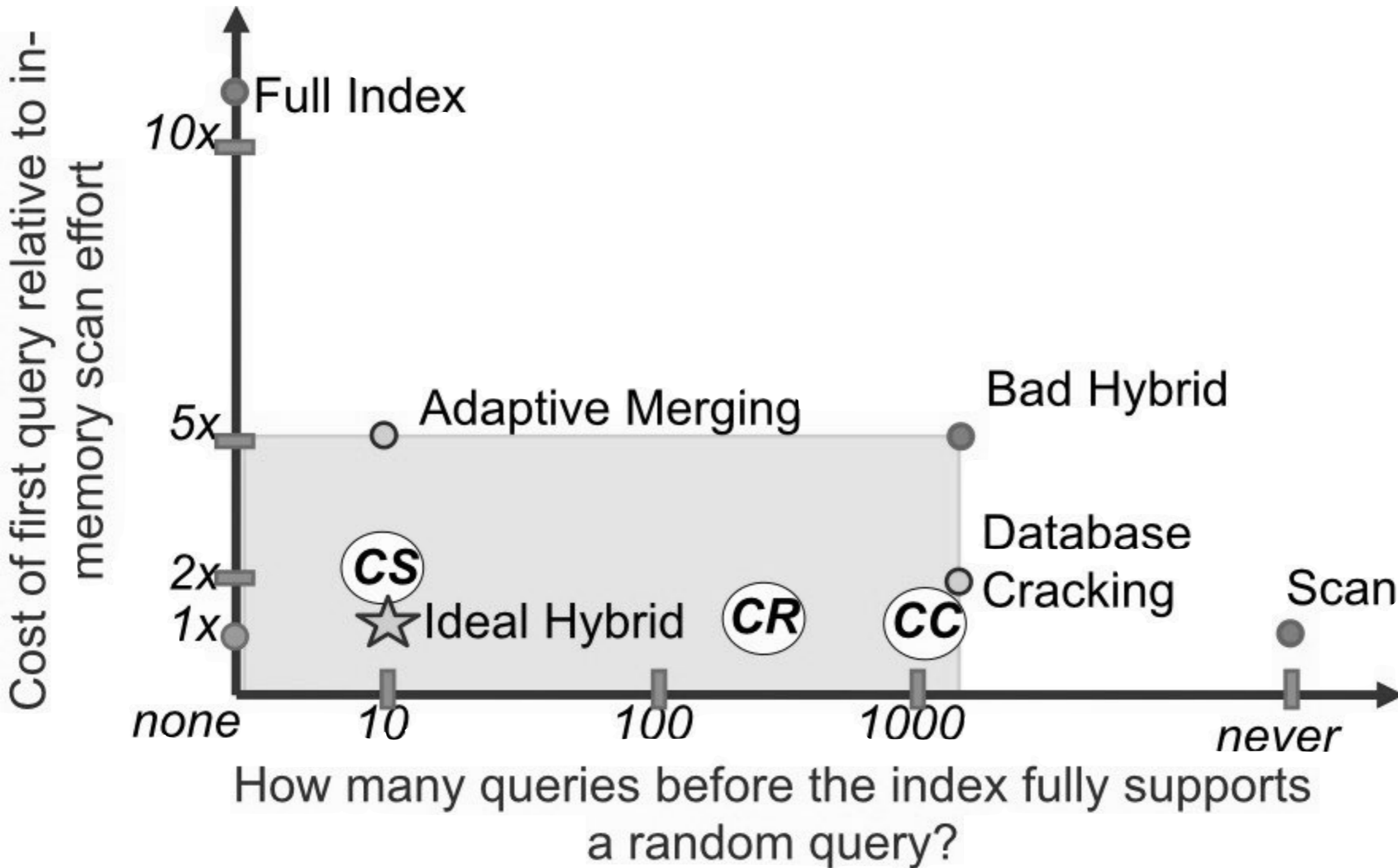




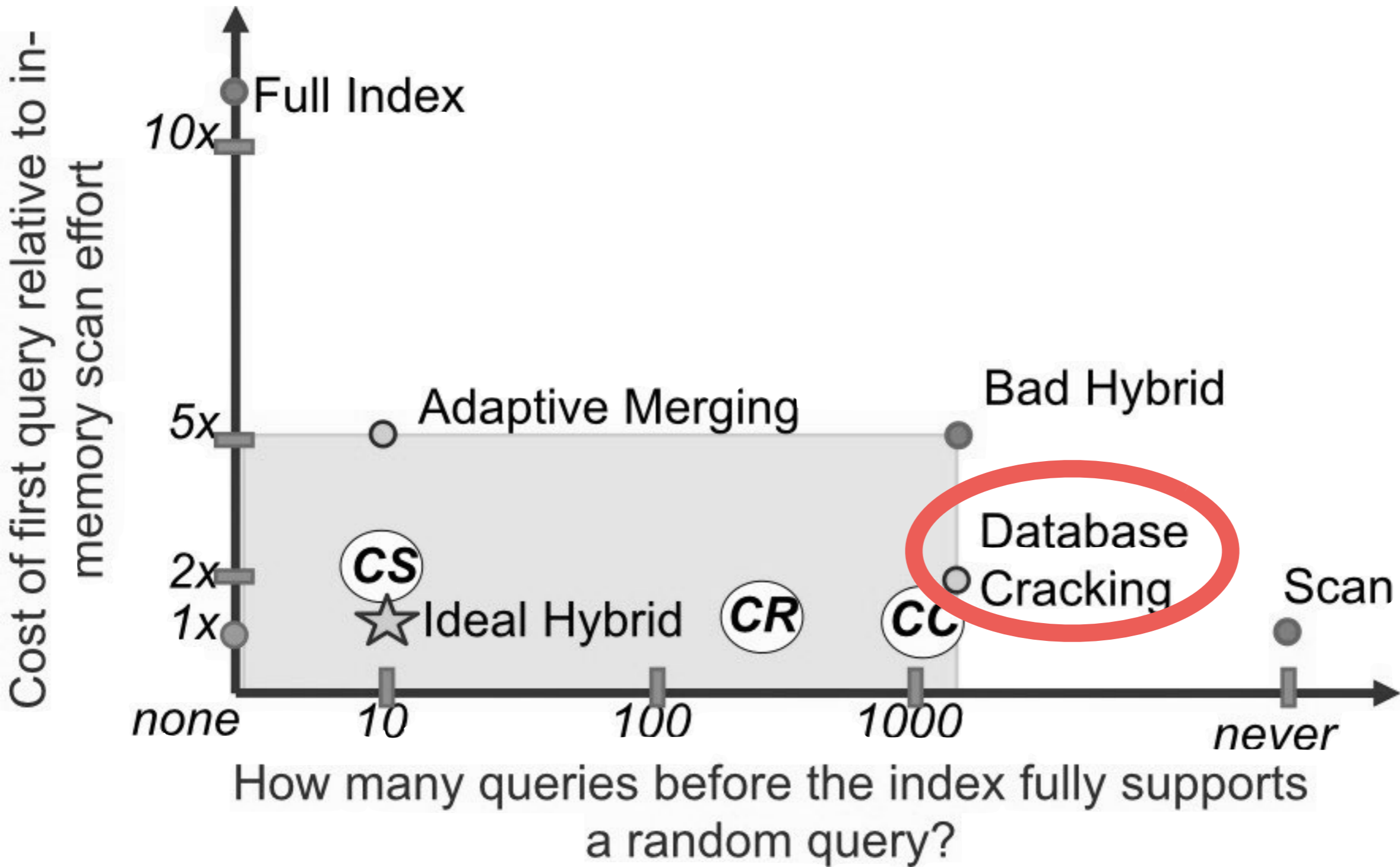




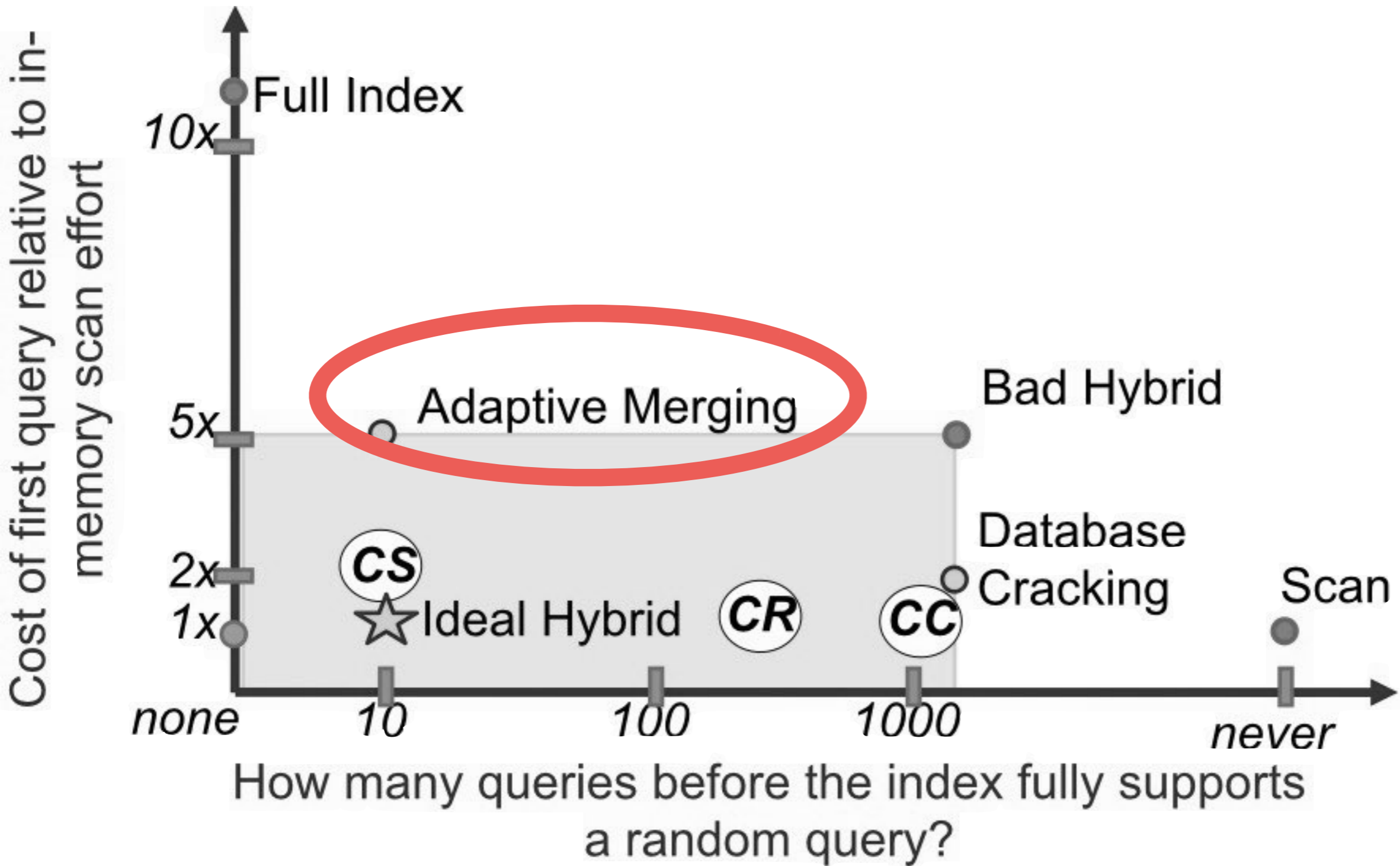
initialization / convergence tradeoff



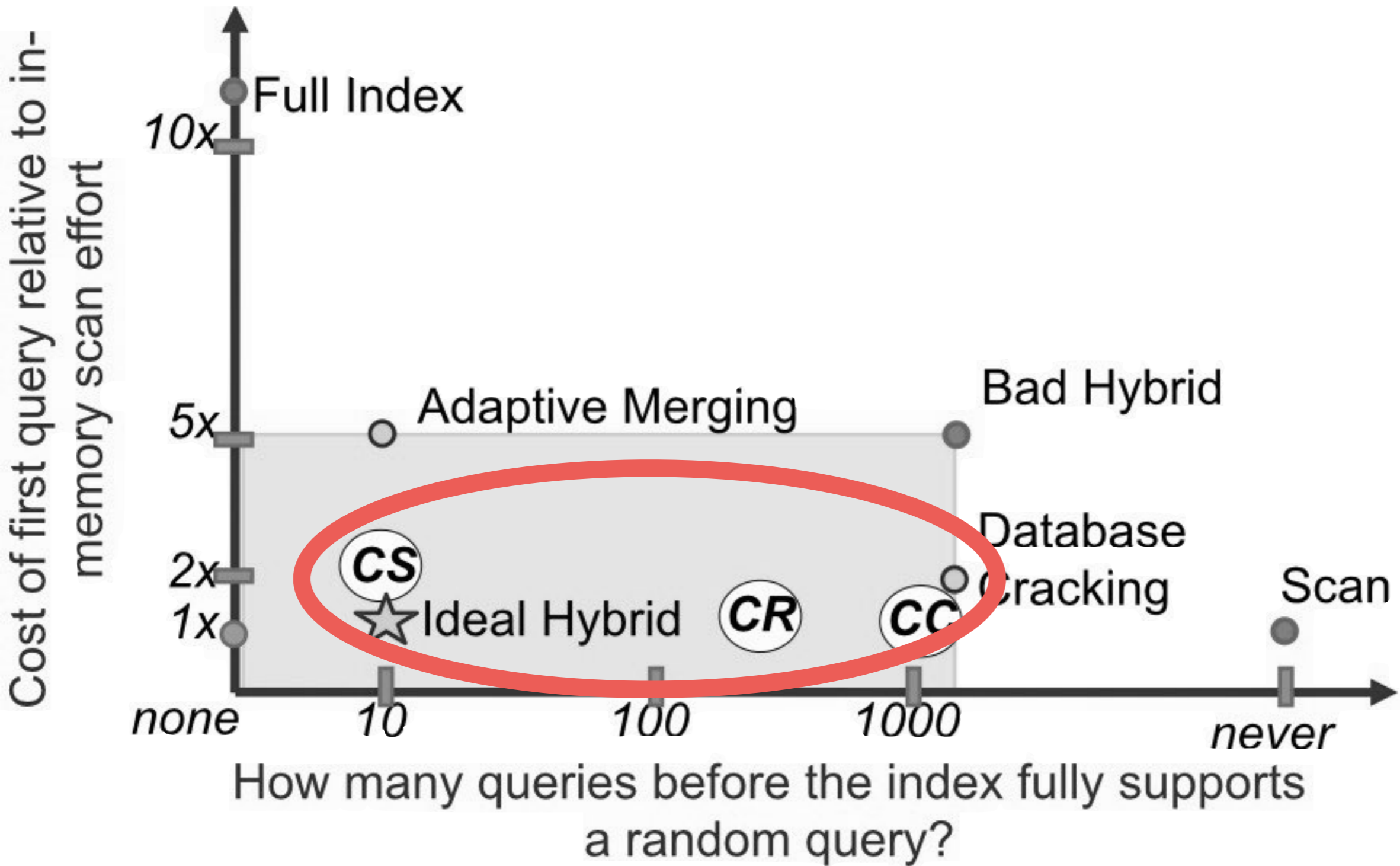
initialization / convergence tradeoff

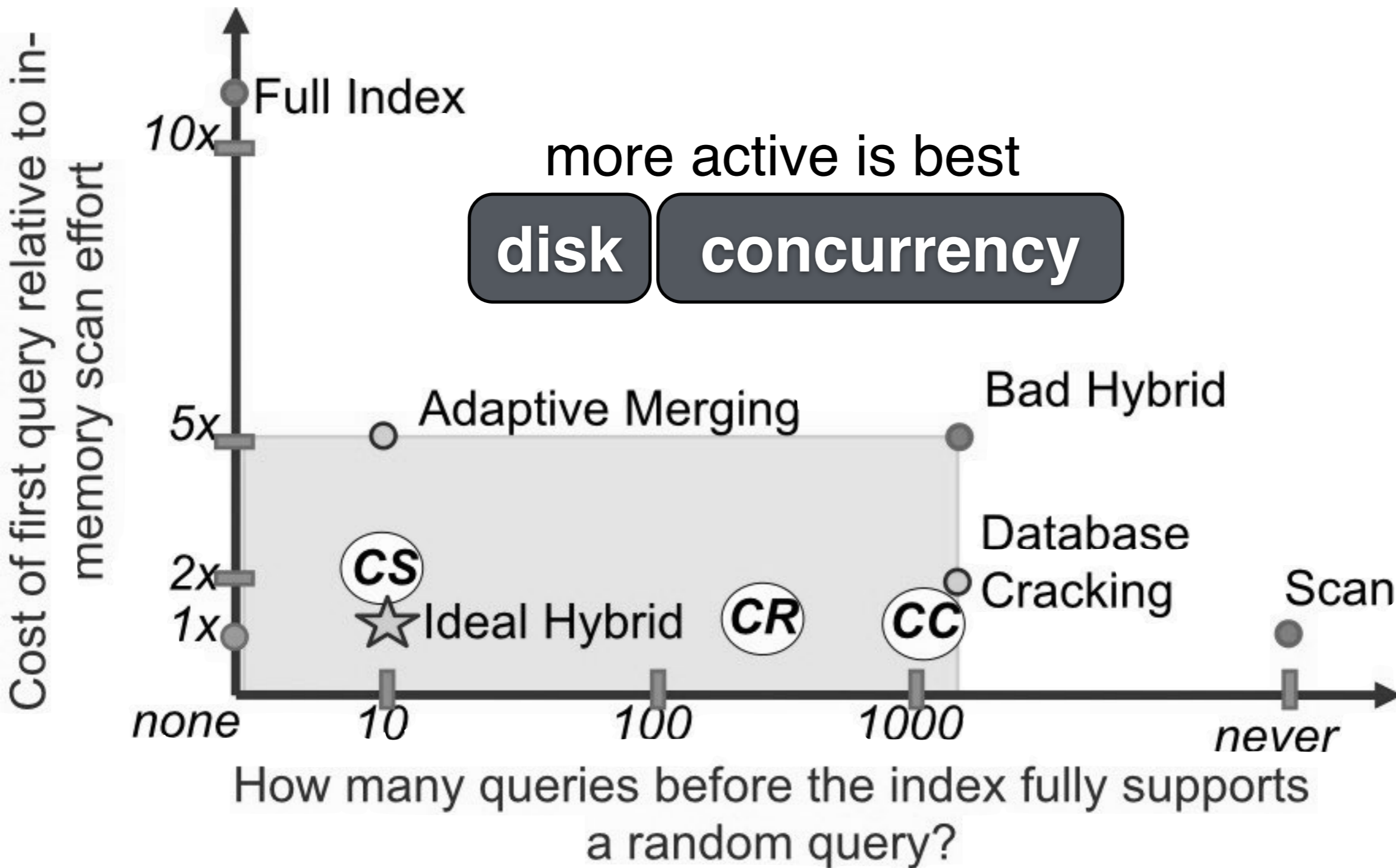


initialization / convergence tradeoff

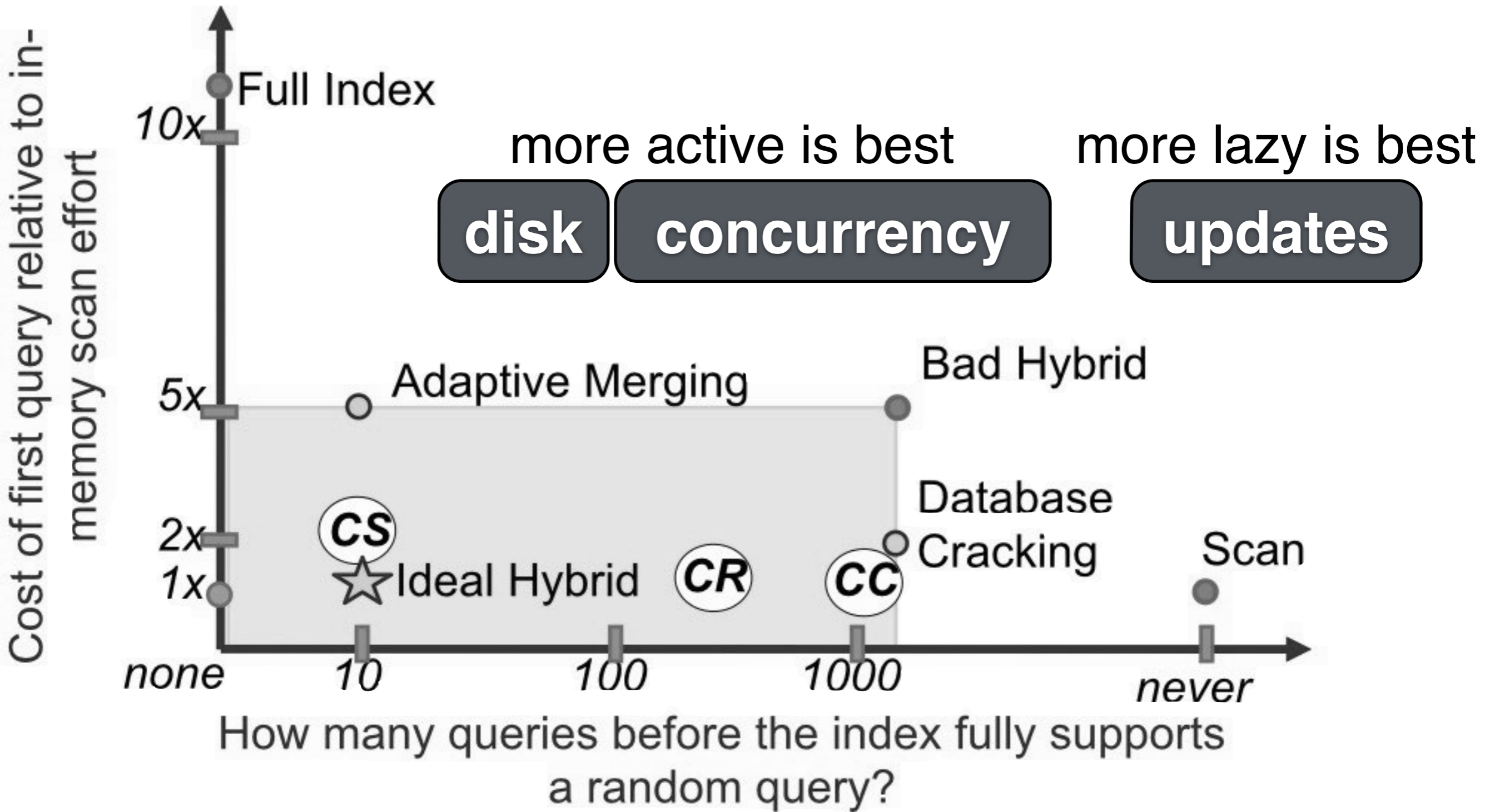


initialization / convergence tradeoff



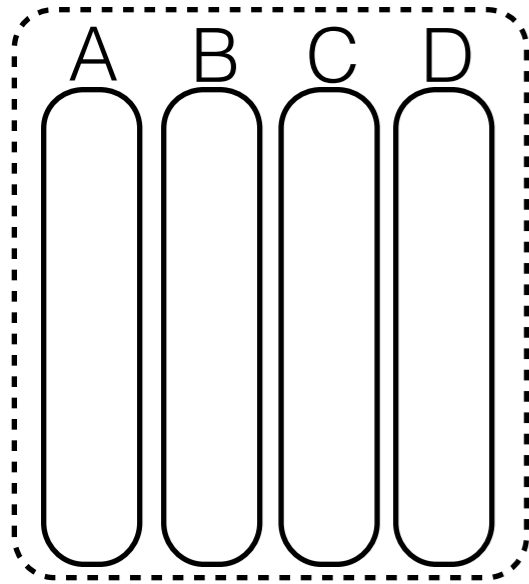


initialization / convergence tradeoff



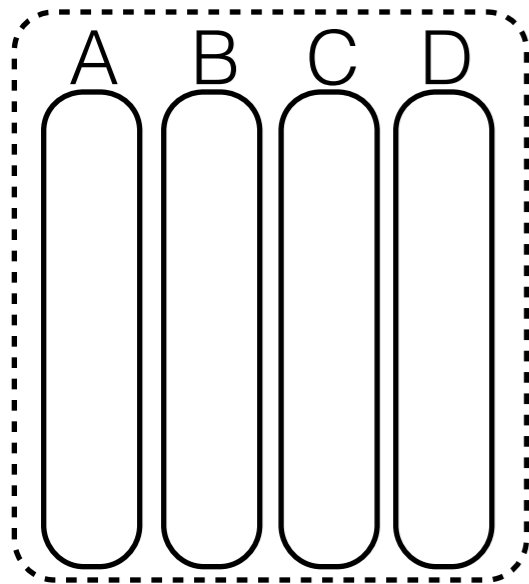
cracking tangram

base data
table 1



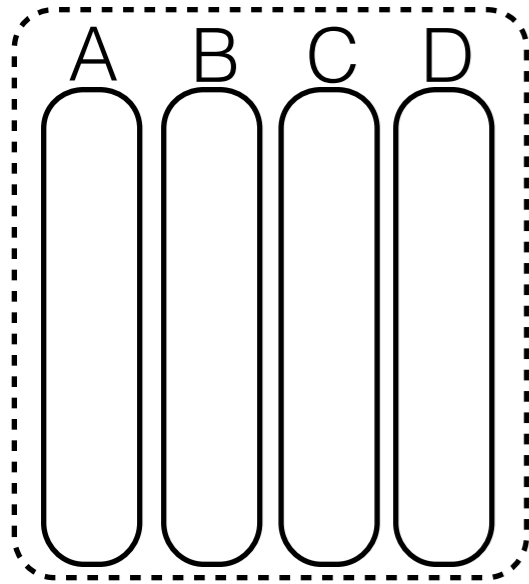
as queries arrive...

table 2



cracking tangram

base data
table 1



as queries arrive...
table 1

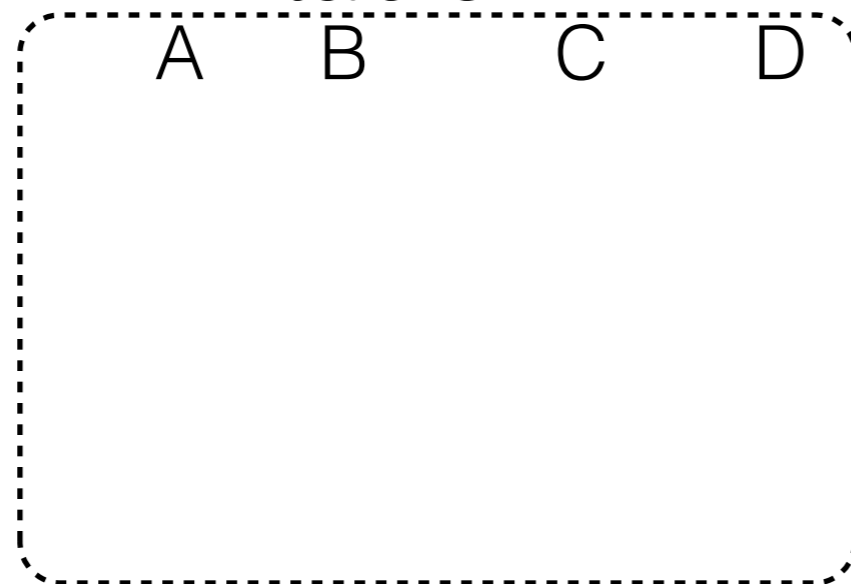


table 2

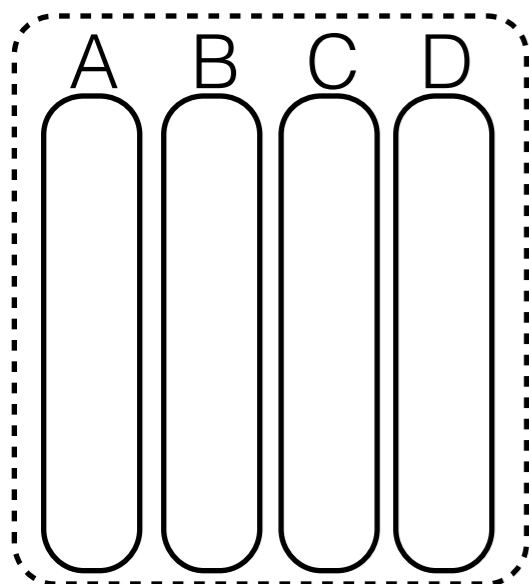
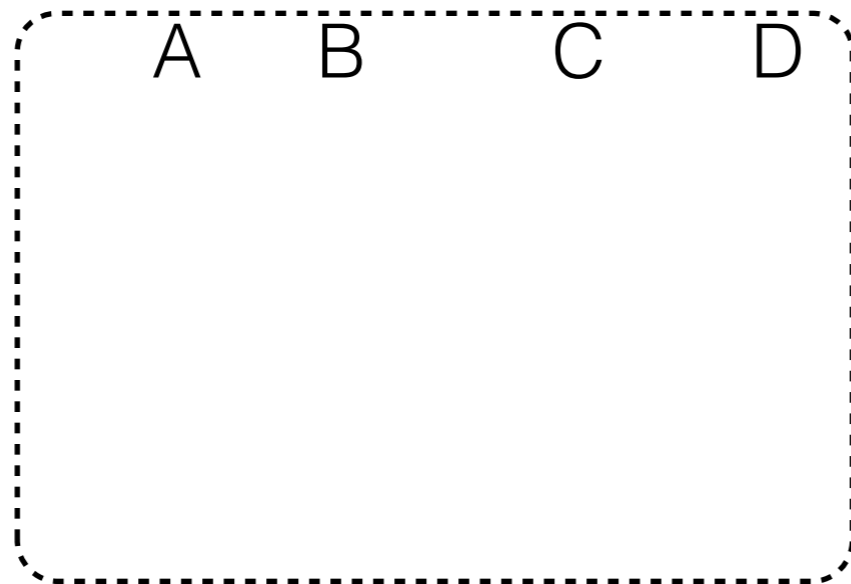
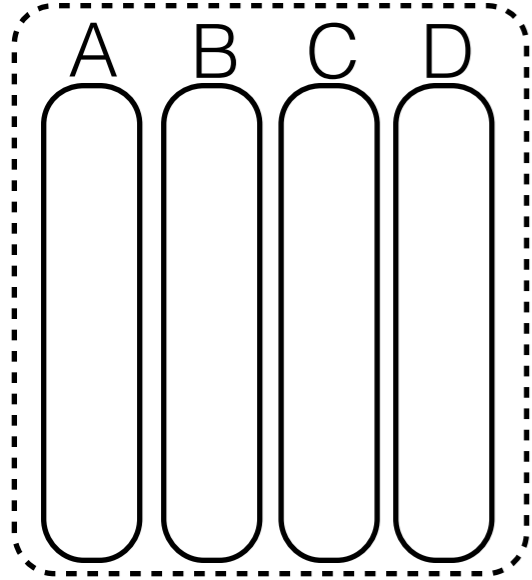


table 2

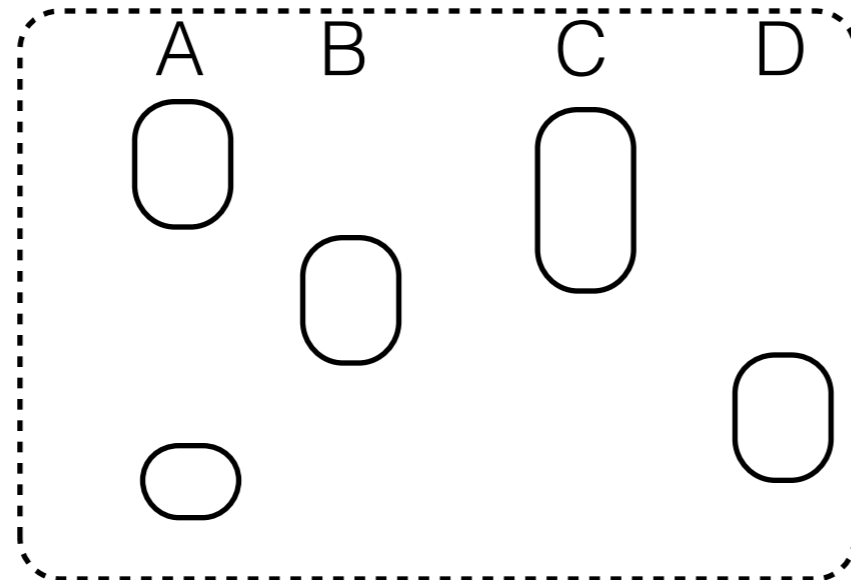


cracking tangram

base data
table 1



as queries arrive...
table 1



partial materialization

table 2

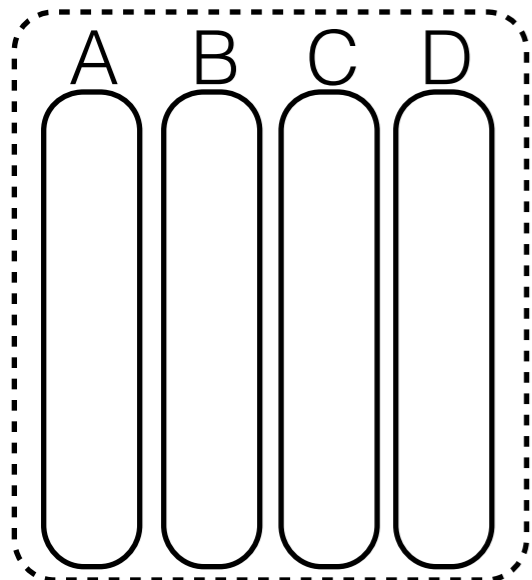
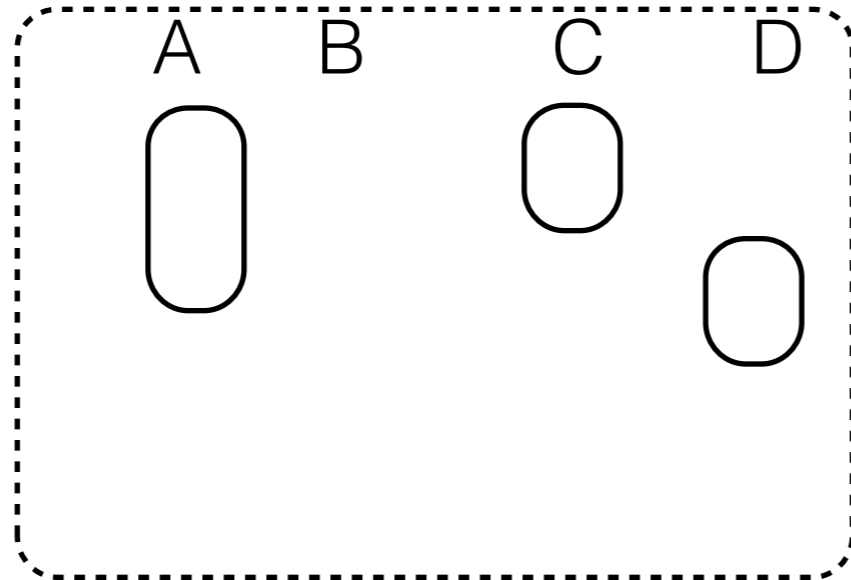
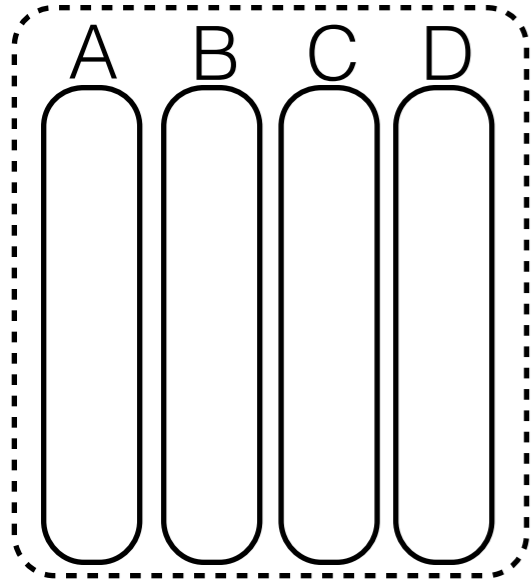


table 2

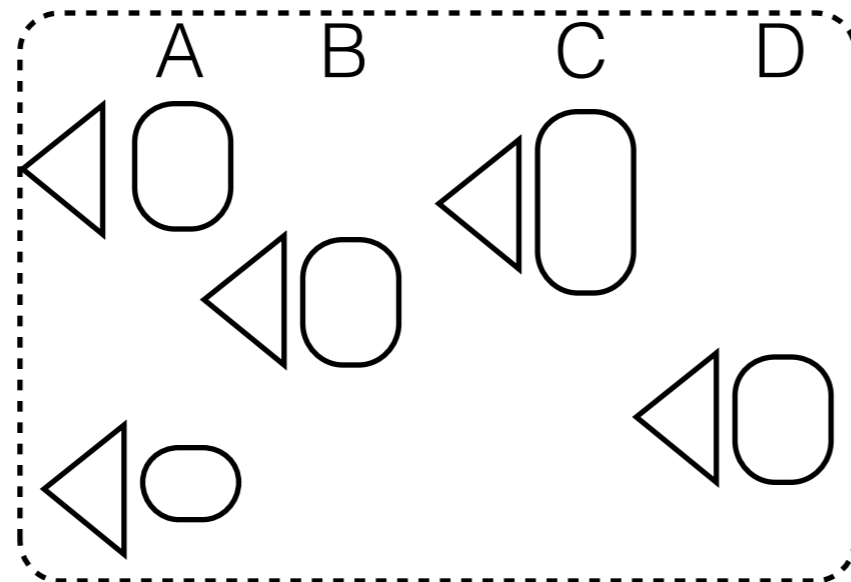


cracking tangram

base data
table 1



as queries arrive...
table 1



partial materialization
partial indexing

table 2

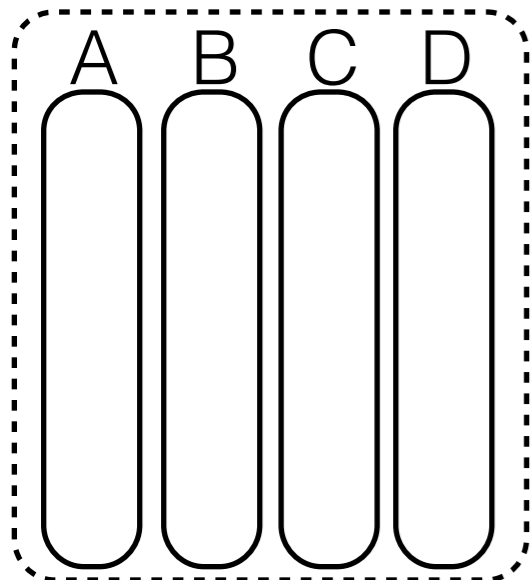
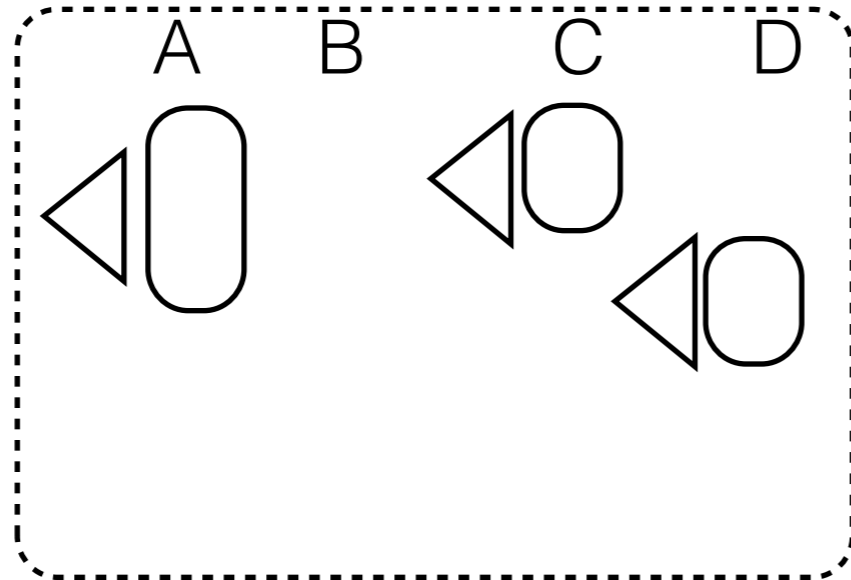
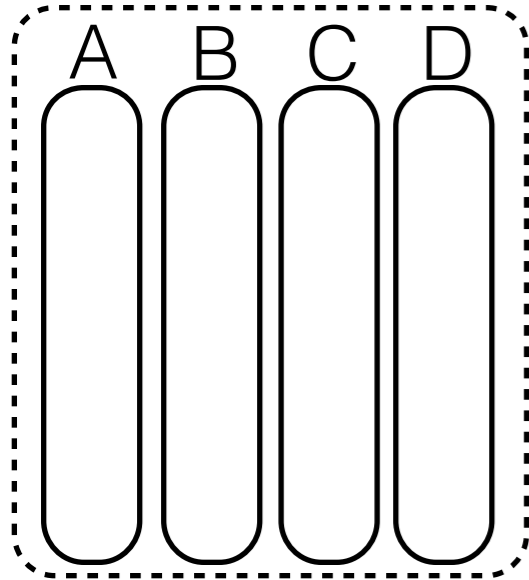


table 2

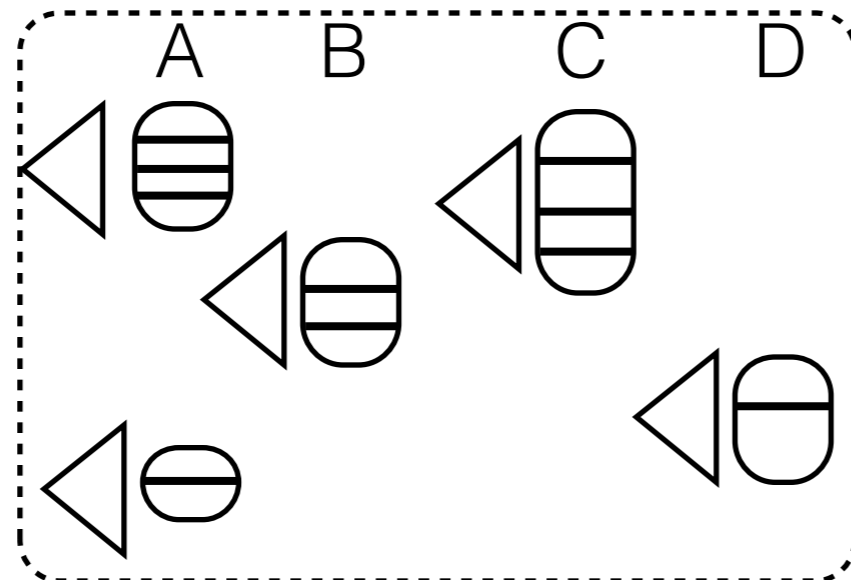


cracking tangram

base data
table 1



as queries arrive...
table 1



partial materialization
partial indexing
continuous adaptation

table 2

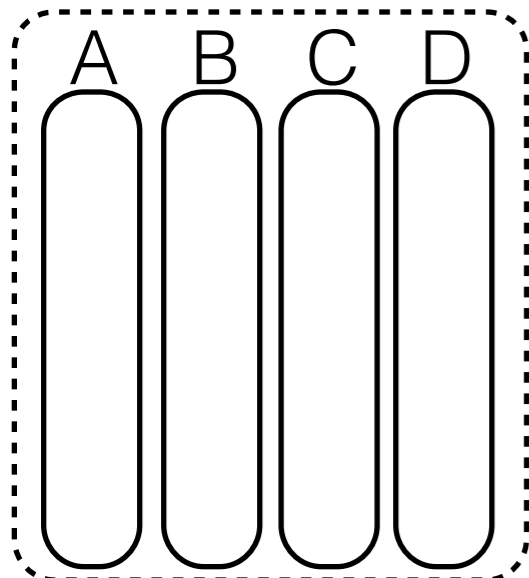
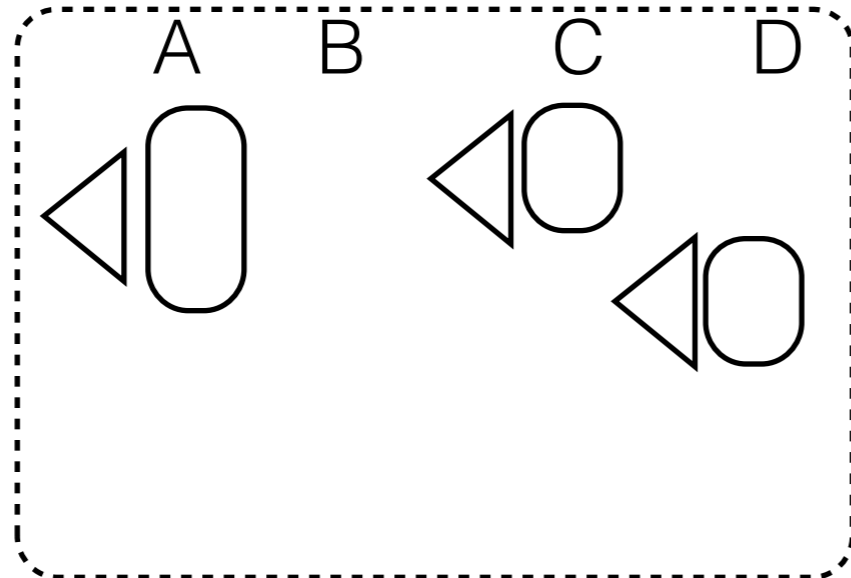
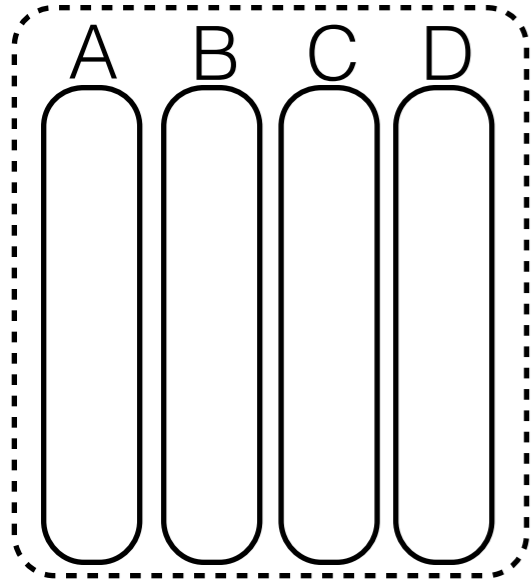


table 2

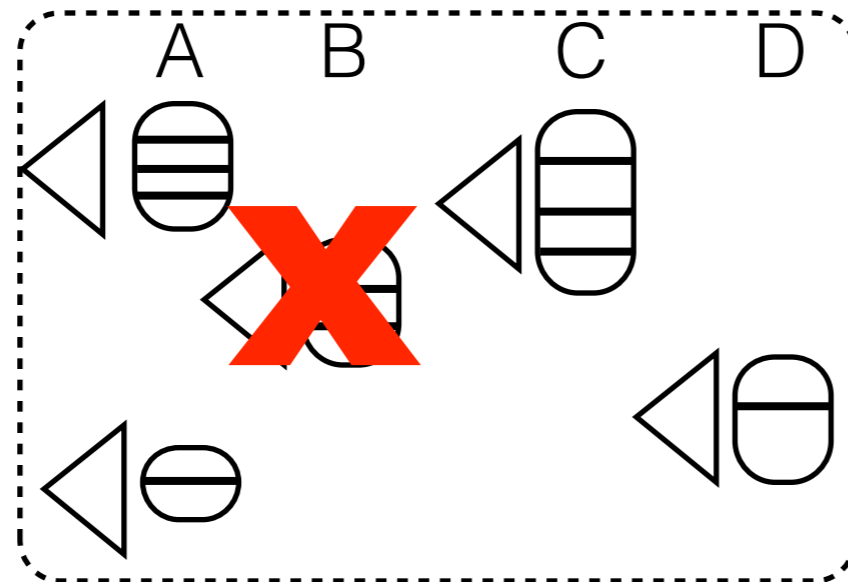


cracking tangram

base data
table 1



as queries arrive...
table 1



- partial materialization
- partial indexing
- continuous adaptation
- storage adaptation

table 2

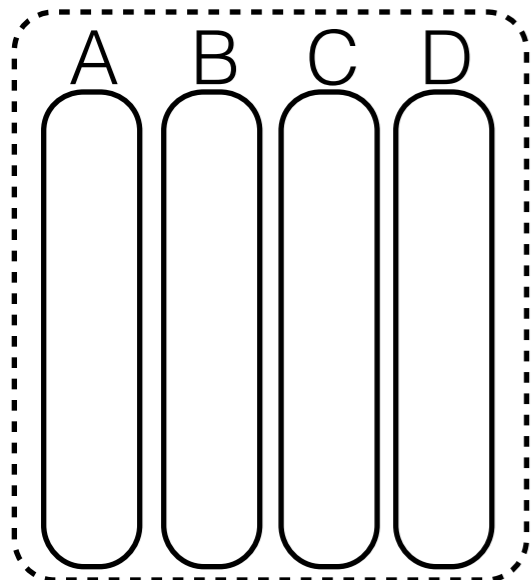
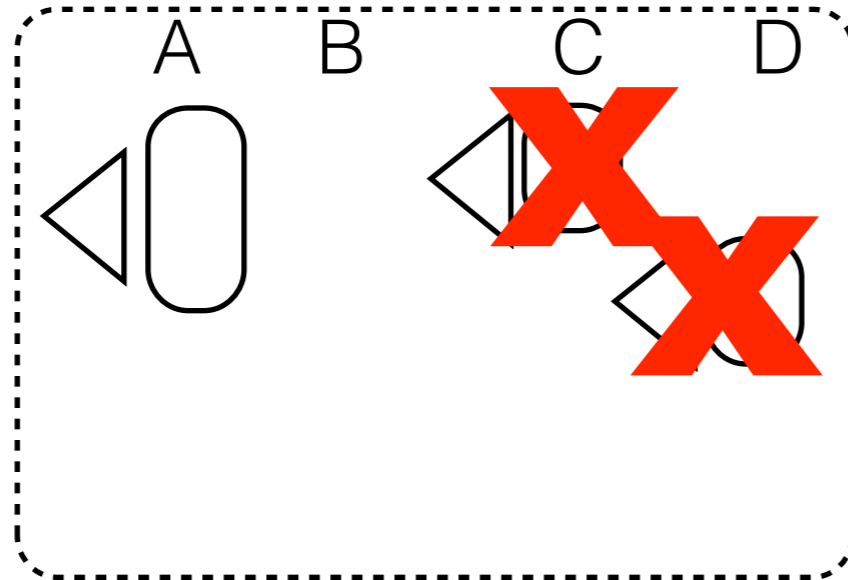
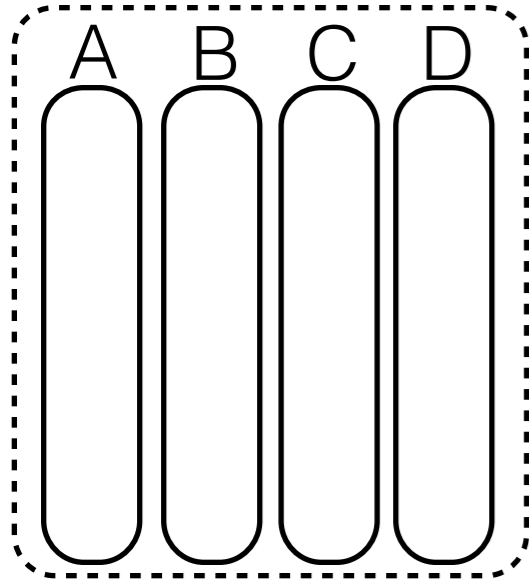


table 2

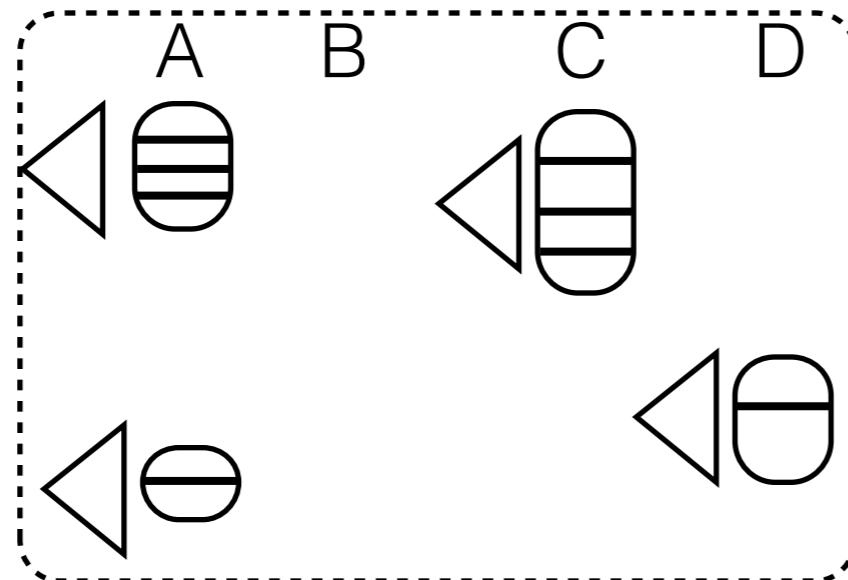


cracking tangram

base data
table 1



as queries arrive...
table 1



partial materialization
partial indexing
continuous adaptation
storage adaptation

table 2

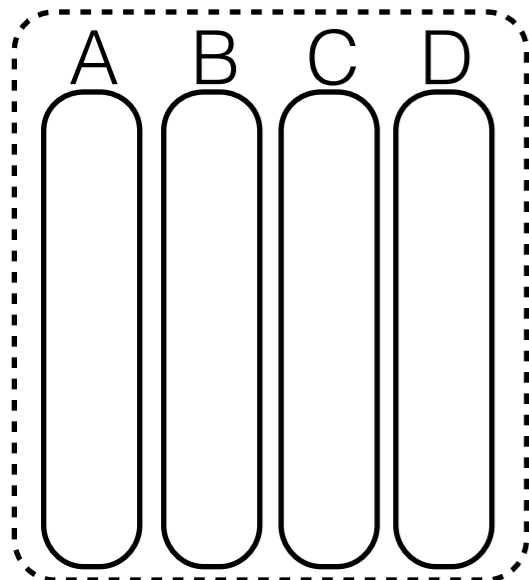
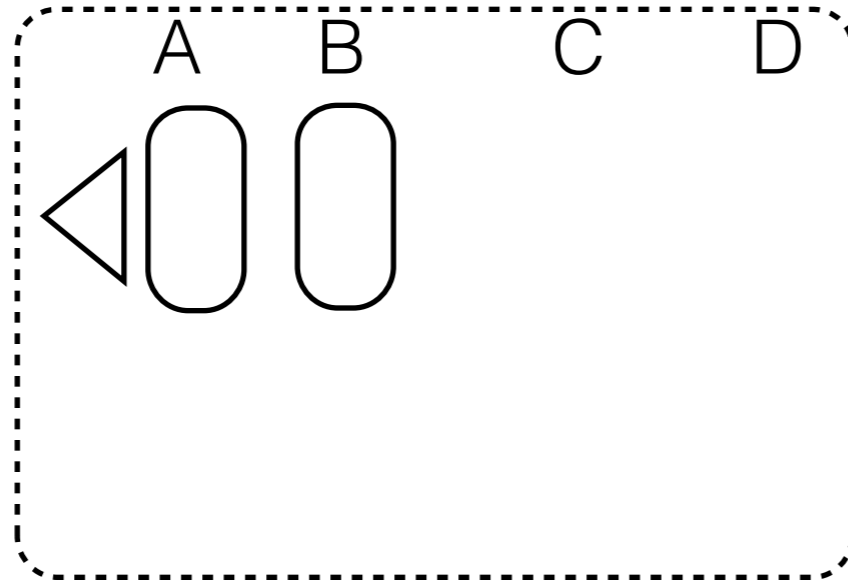
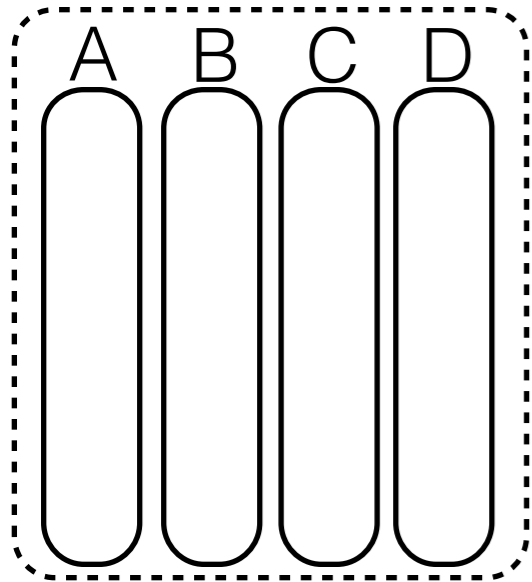


table 2



cracking tangram

base data
table 1



as queries arrive...
table 1

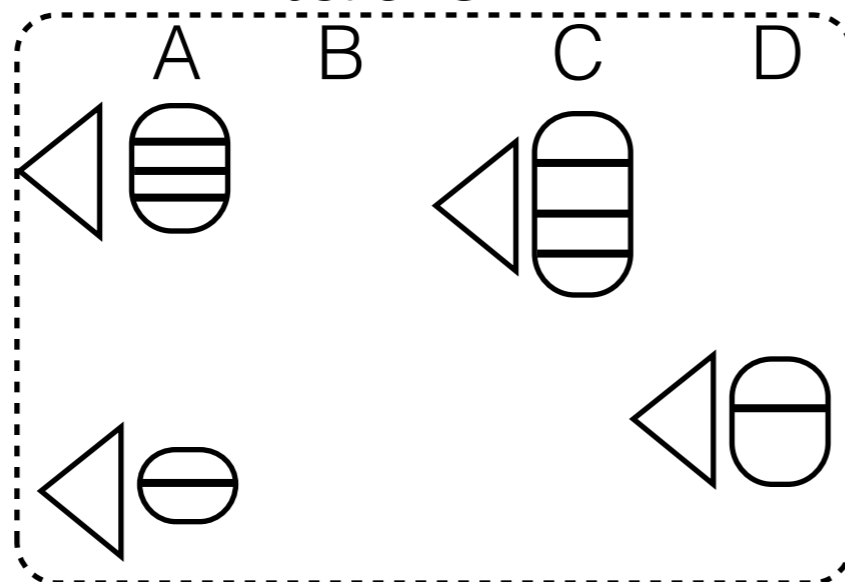


table 2

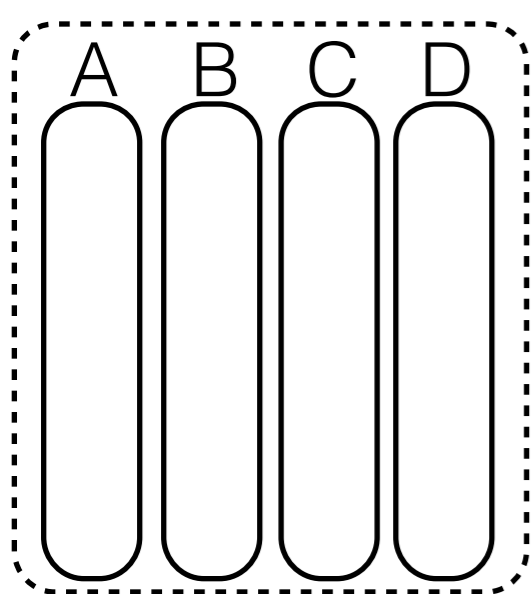
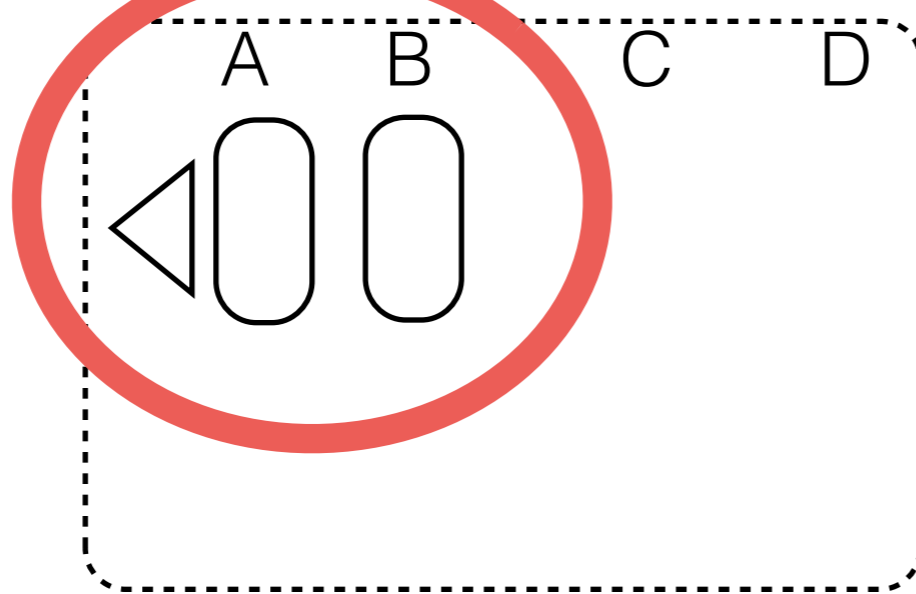


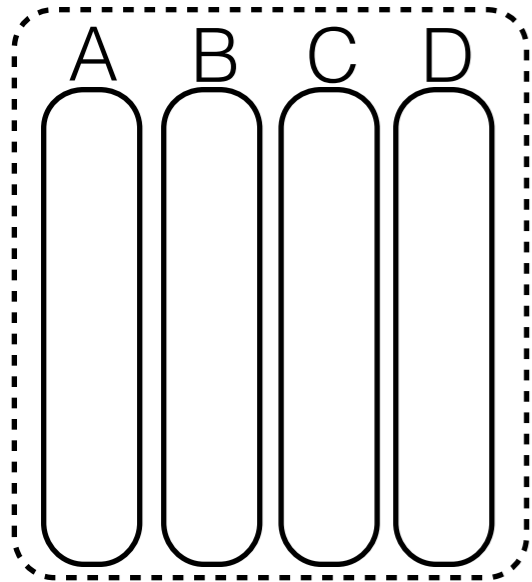
table 2



- partial materialization
- partial indexing
- continuous adaptation
- storage adaptation
- no tuple reconstruction

cracking tangram

base data
table 1



as queries arrive...
table 1

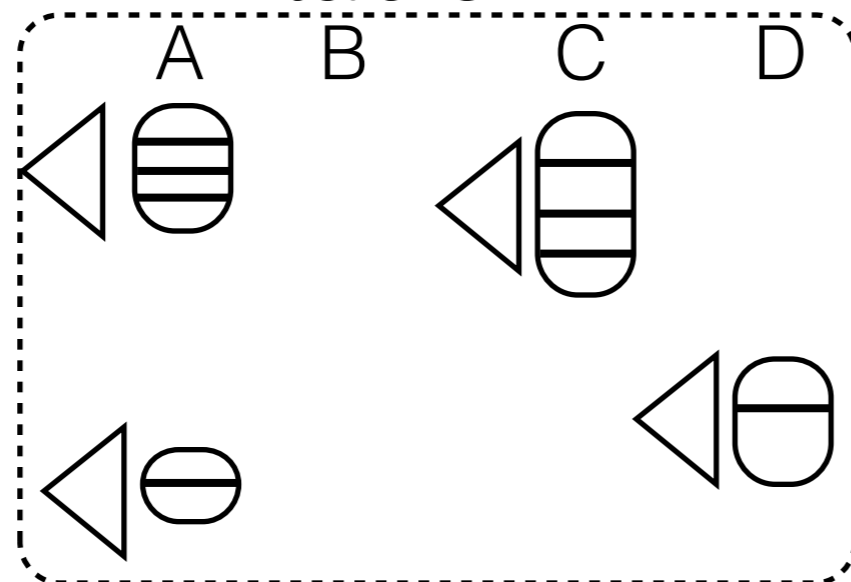


table 2

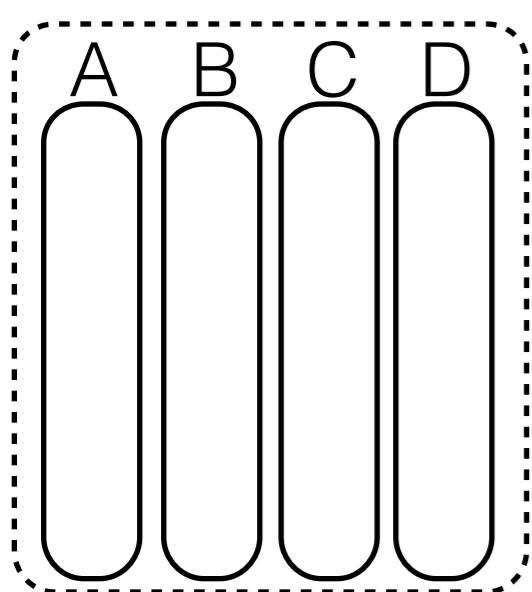
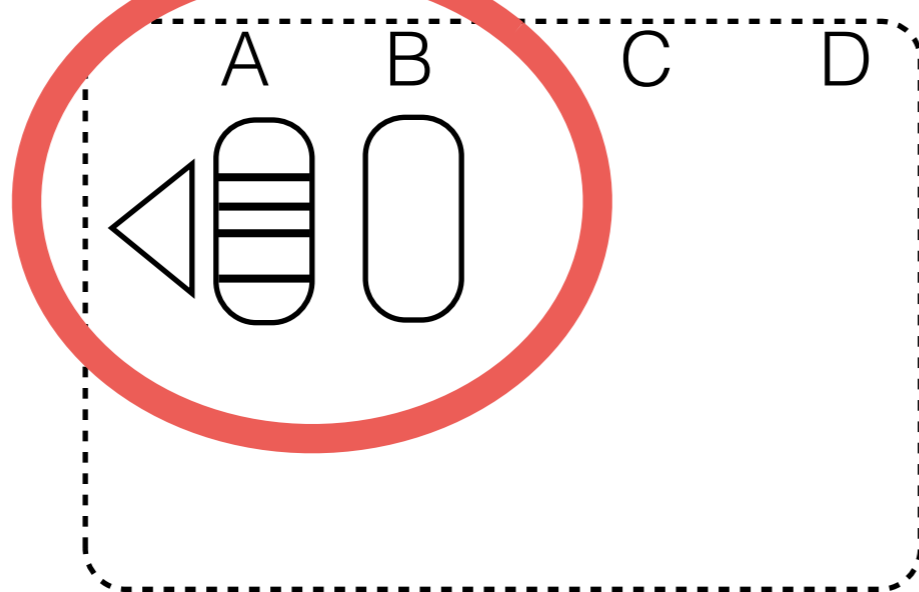


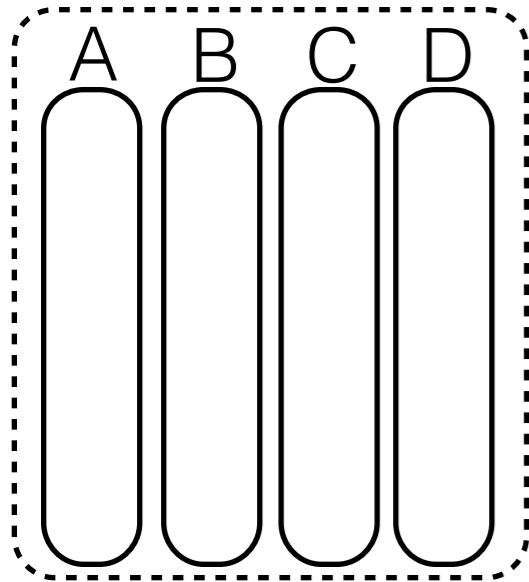
table 2



- partial materialization
- partial indexing
- continuous adaptation
- storage adaptation
- no tuple reconstruction
- adaptive alignment

cracking tangram

base data
table 1



as queries arrive...
table 1

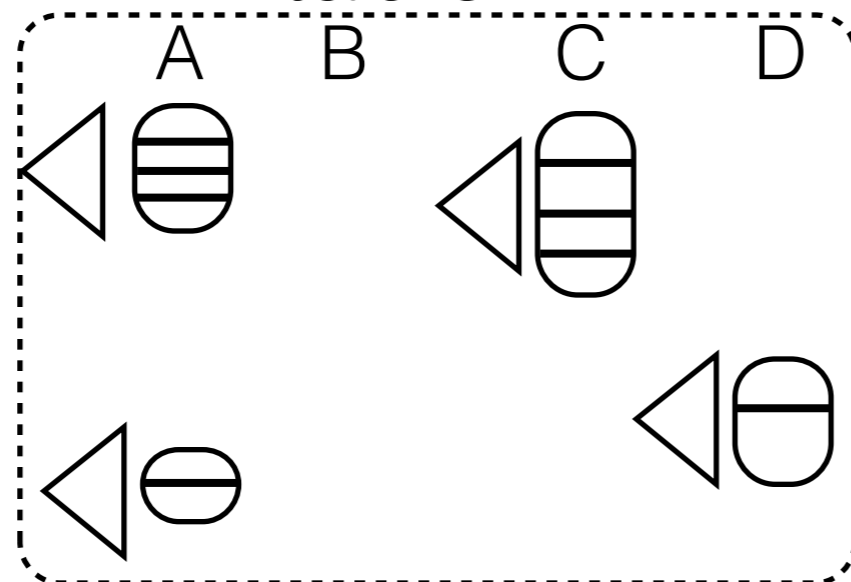


table 2

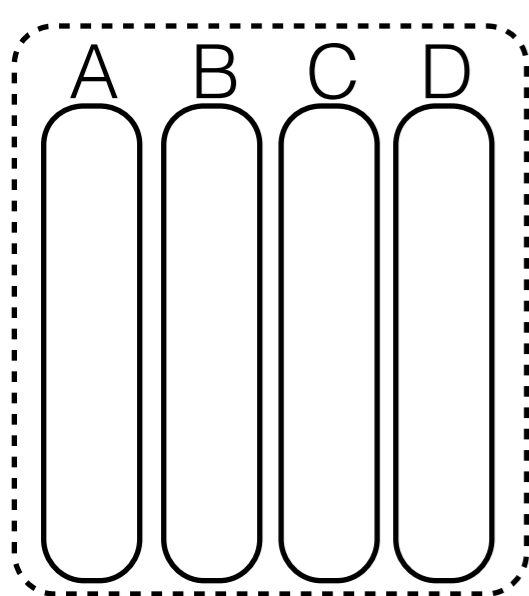
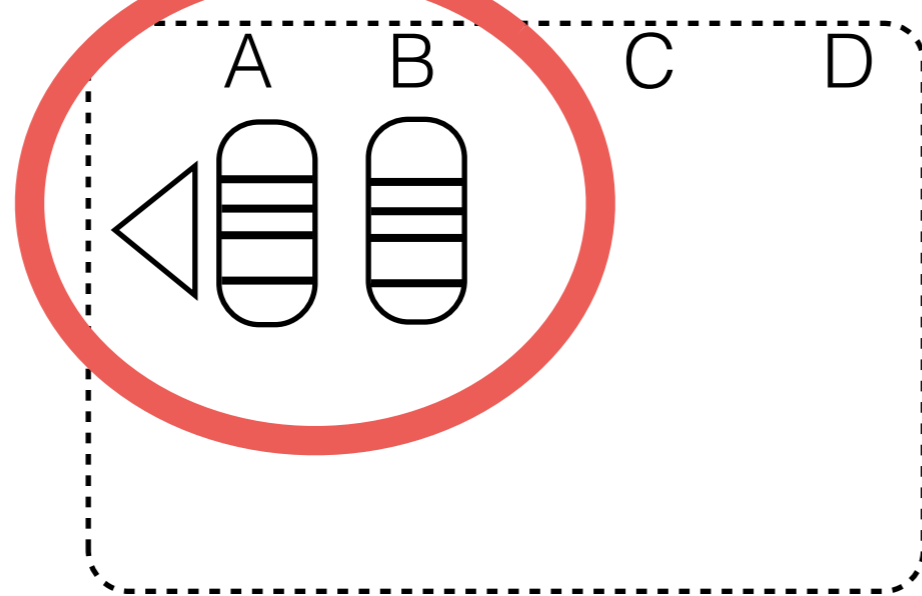


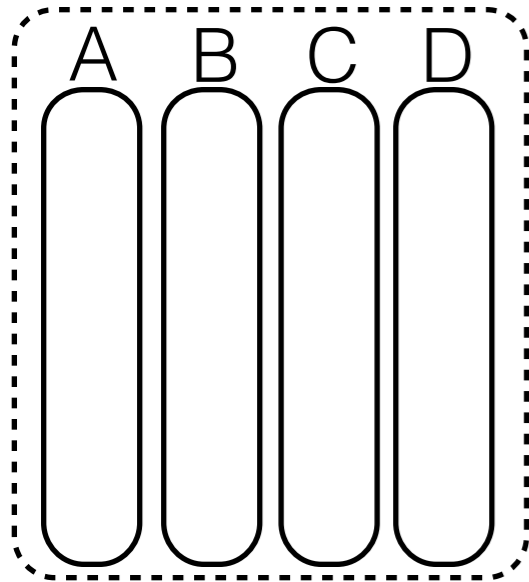
table 2



- partial materialization
- partial indexing
- continuous adaptation
- storage adaptation
- no tuple reconstruction
- adaptive alignment

cracking tangram

base data
table 1



as queries arrive...
table 1

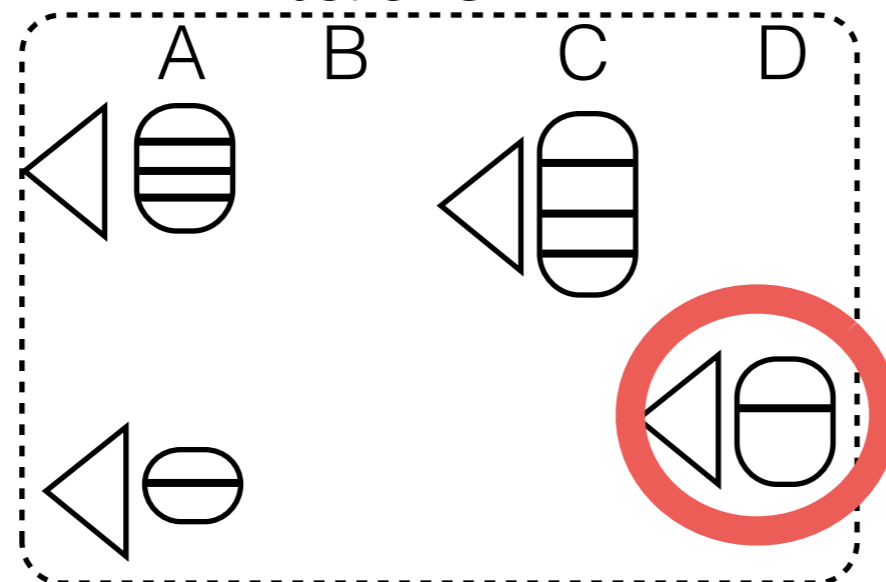


table 2

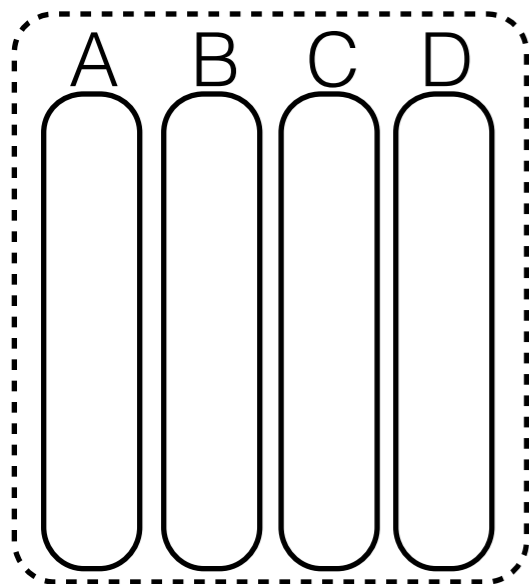
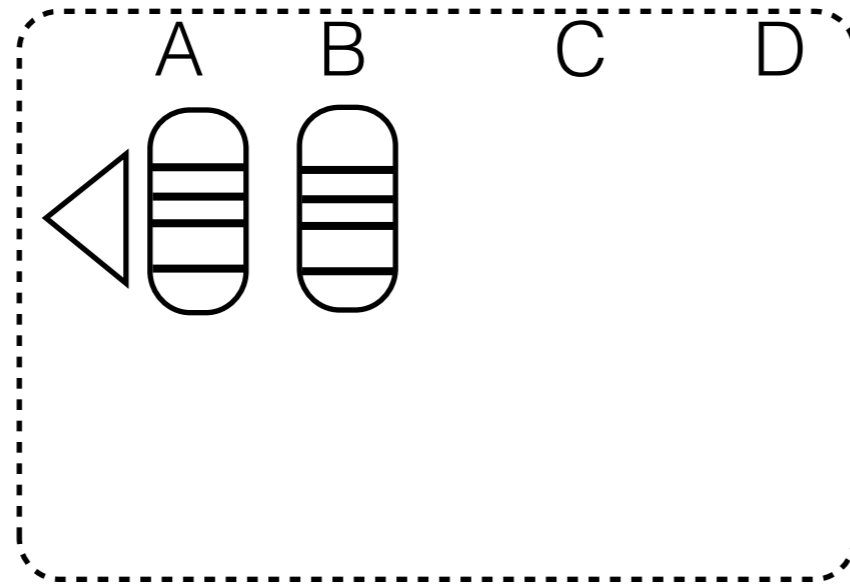


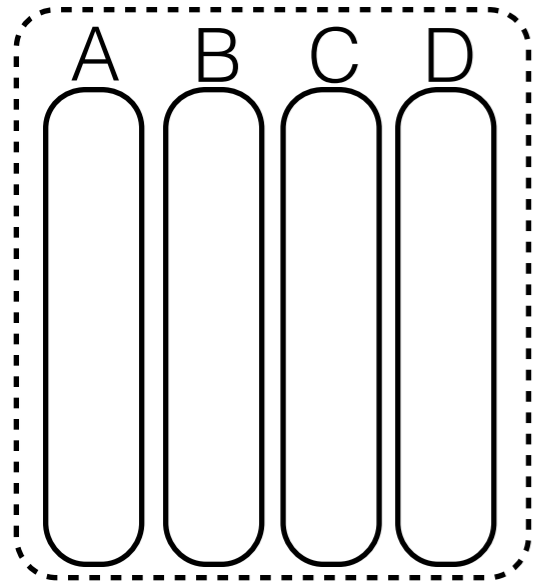
table 2



- partial materialization
- partial indexing
- continuous adaptation
- storage adaptation
- no tuple reconstruction
- adaptive alignment
- sort in caches

cracking tangram

base data
table 1



as queries arrive...
table 1

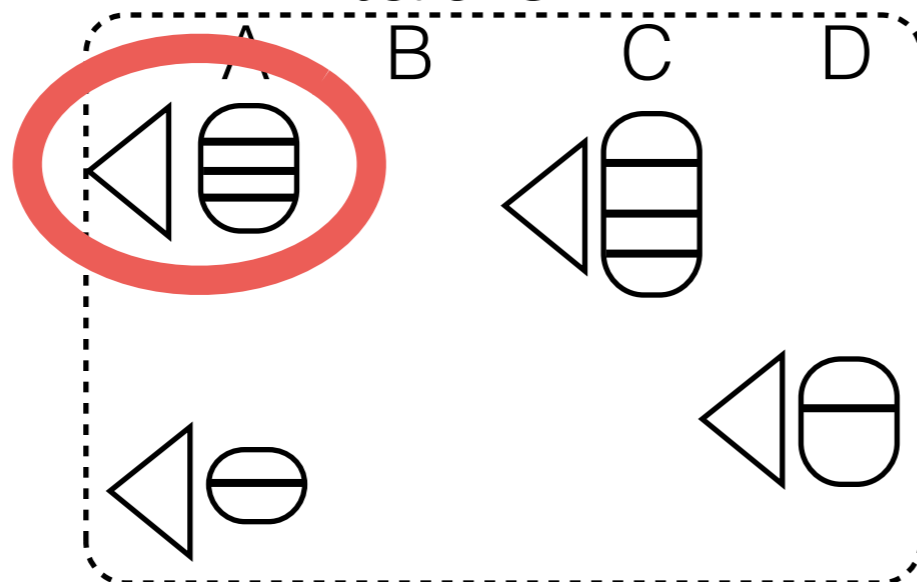


table 2

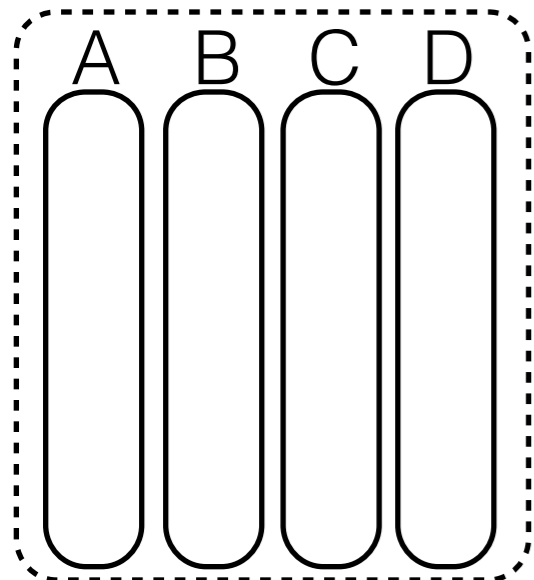
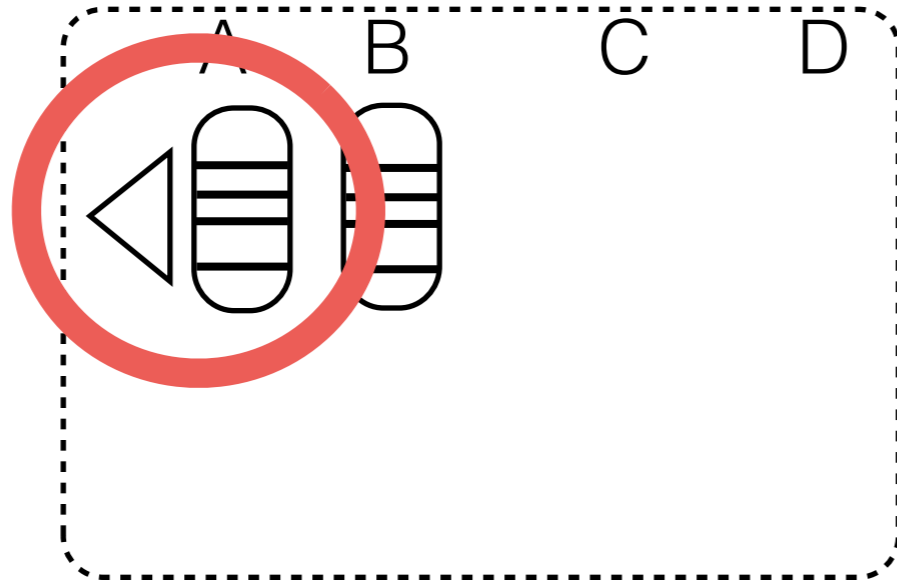


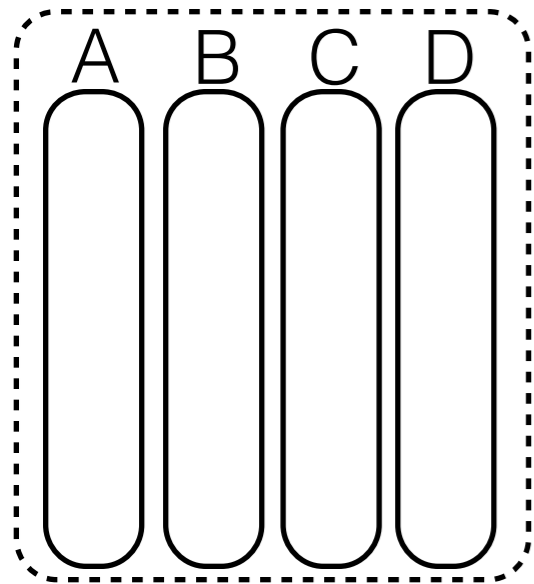
table 2



- partial materialization
- partial indexing
- continuous adaptation
- storage adaptation
- no tuple reconstruction
- adaptive alignment
- sort in caches
- crack joins

cracking tangram

base data
table 1



as queries arrive...
table 1

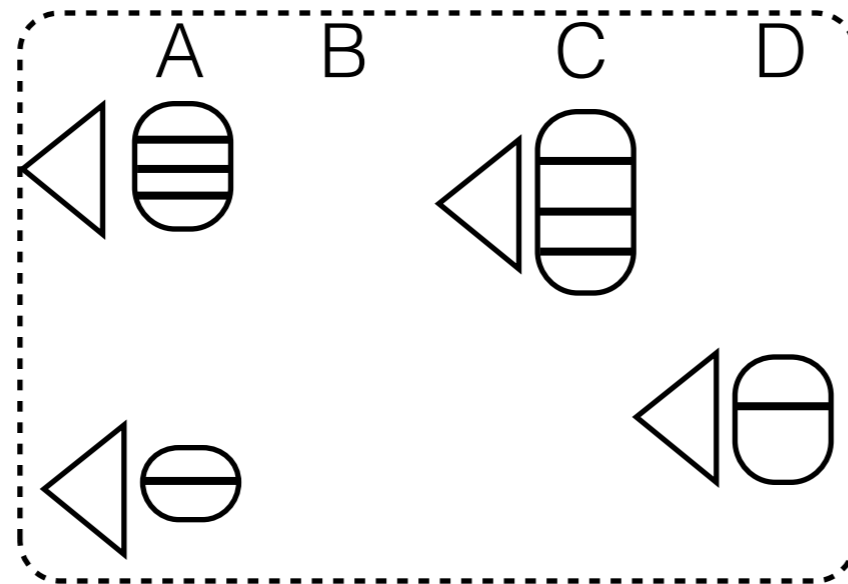


table 2

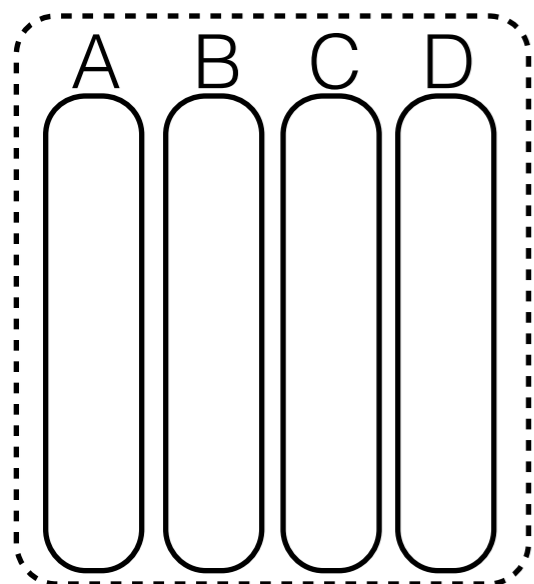
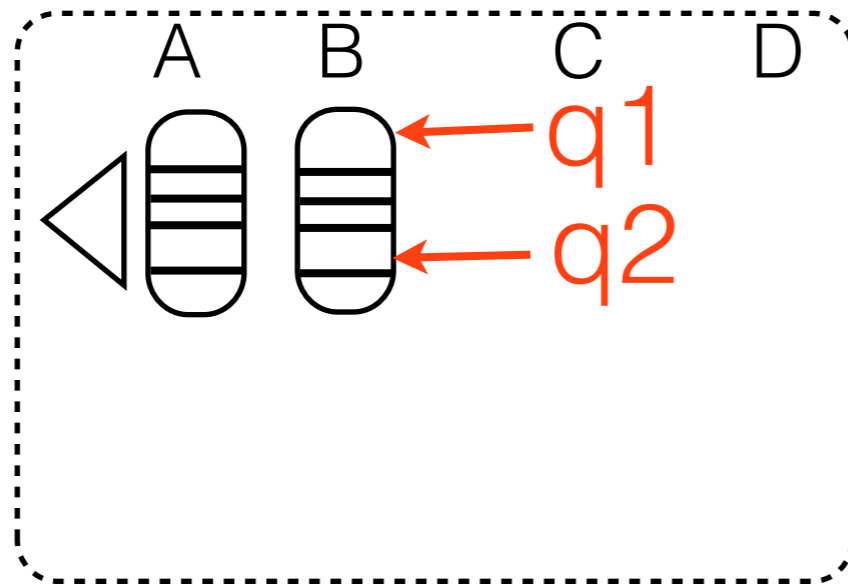


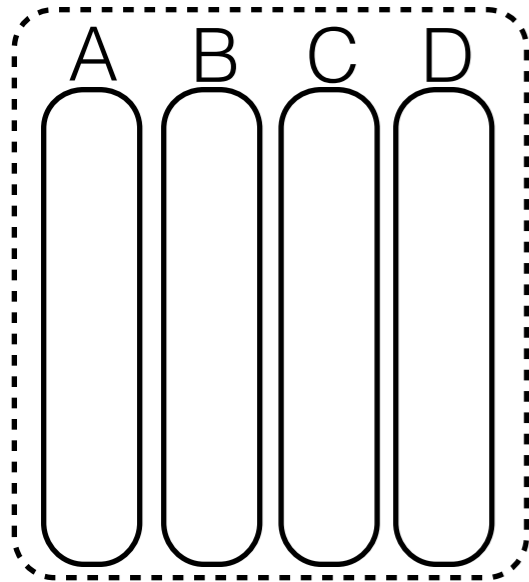
table 2



- partial materialization
- partial indexing
- continuous adaptation
- storage adaptation
- no tuple reconstruction
- adaptive alignment
- sort in caches
- crack joins
- lightweight locking

cracking tangram

base data
table 1



as queries arrive...
table 1

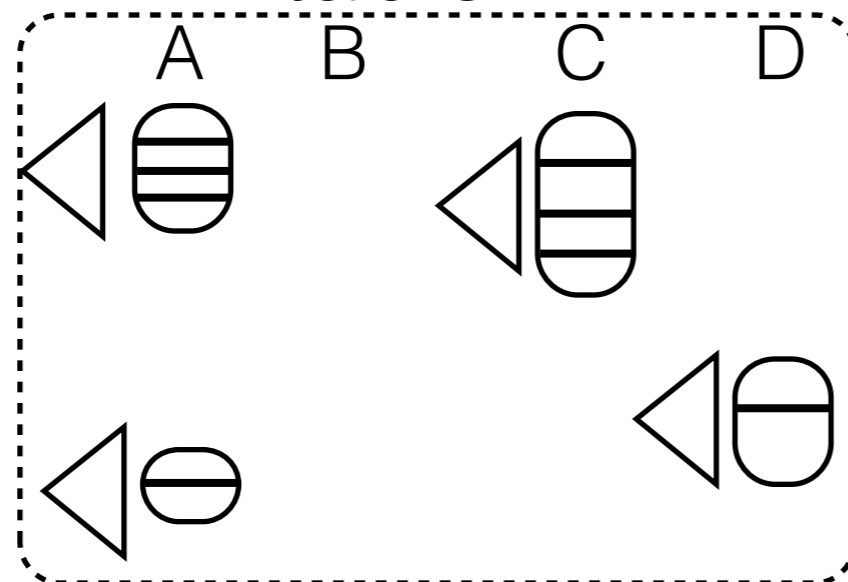


table 2

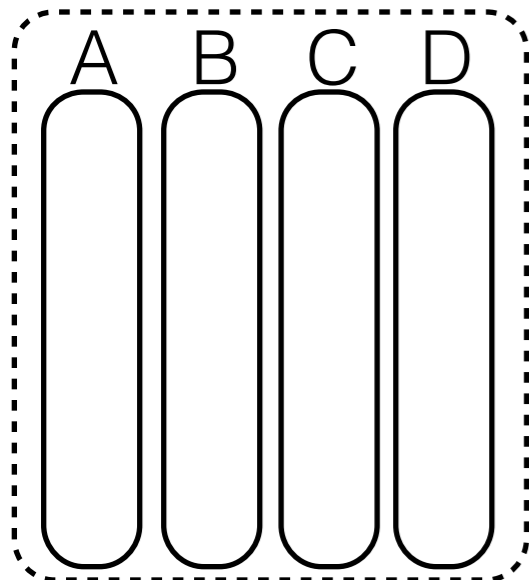
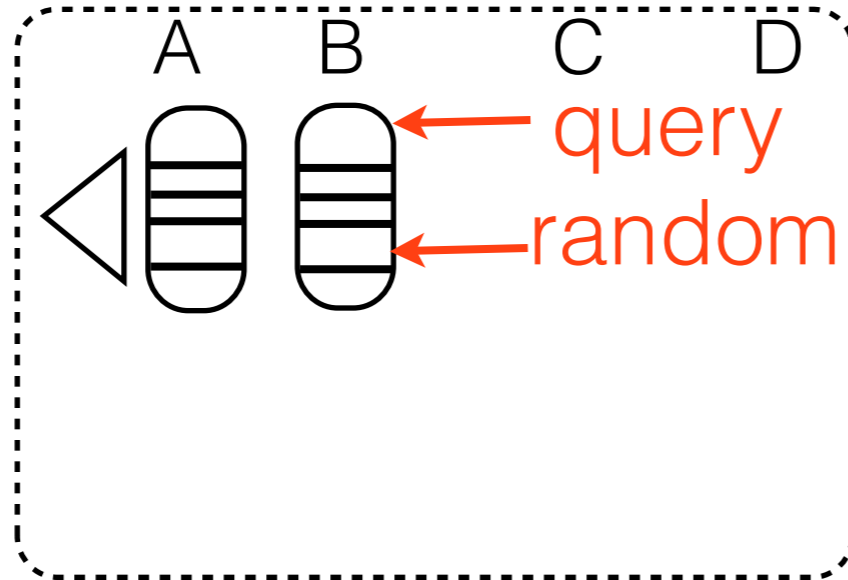


table 2

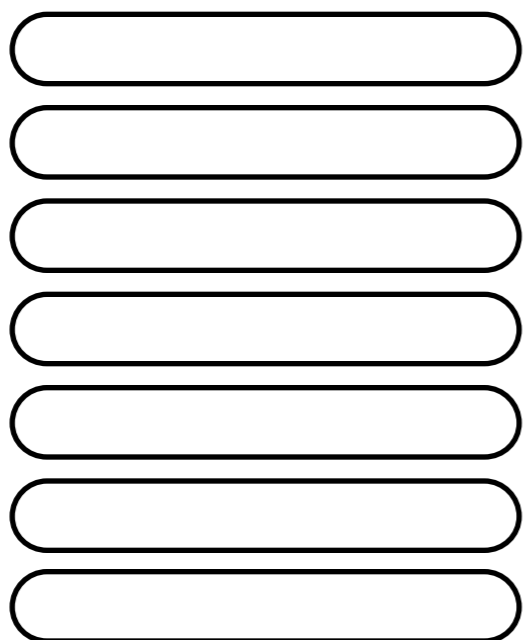


- partial materialization
- partial indexing
- continuous adaptation
- storage adaptation
- no tuple reconstruction
- adaptive alignment
- sort in caches
- crack joins
- lightweight locking
- stochastic cracking

adaptive storage

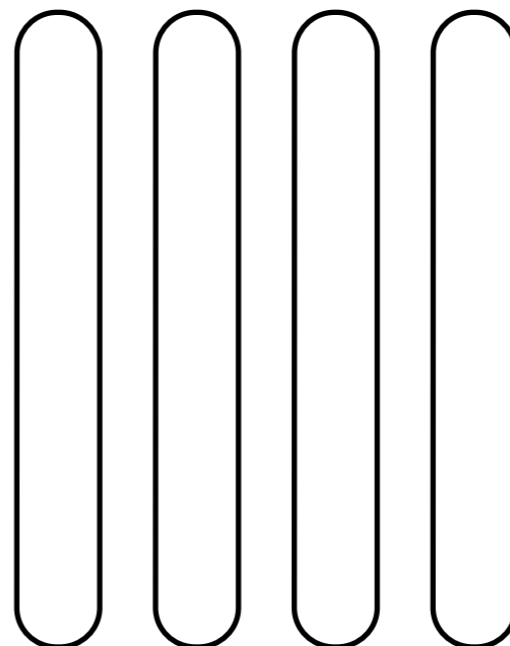
row-store

A B C D

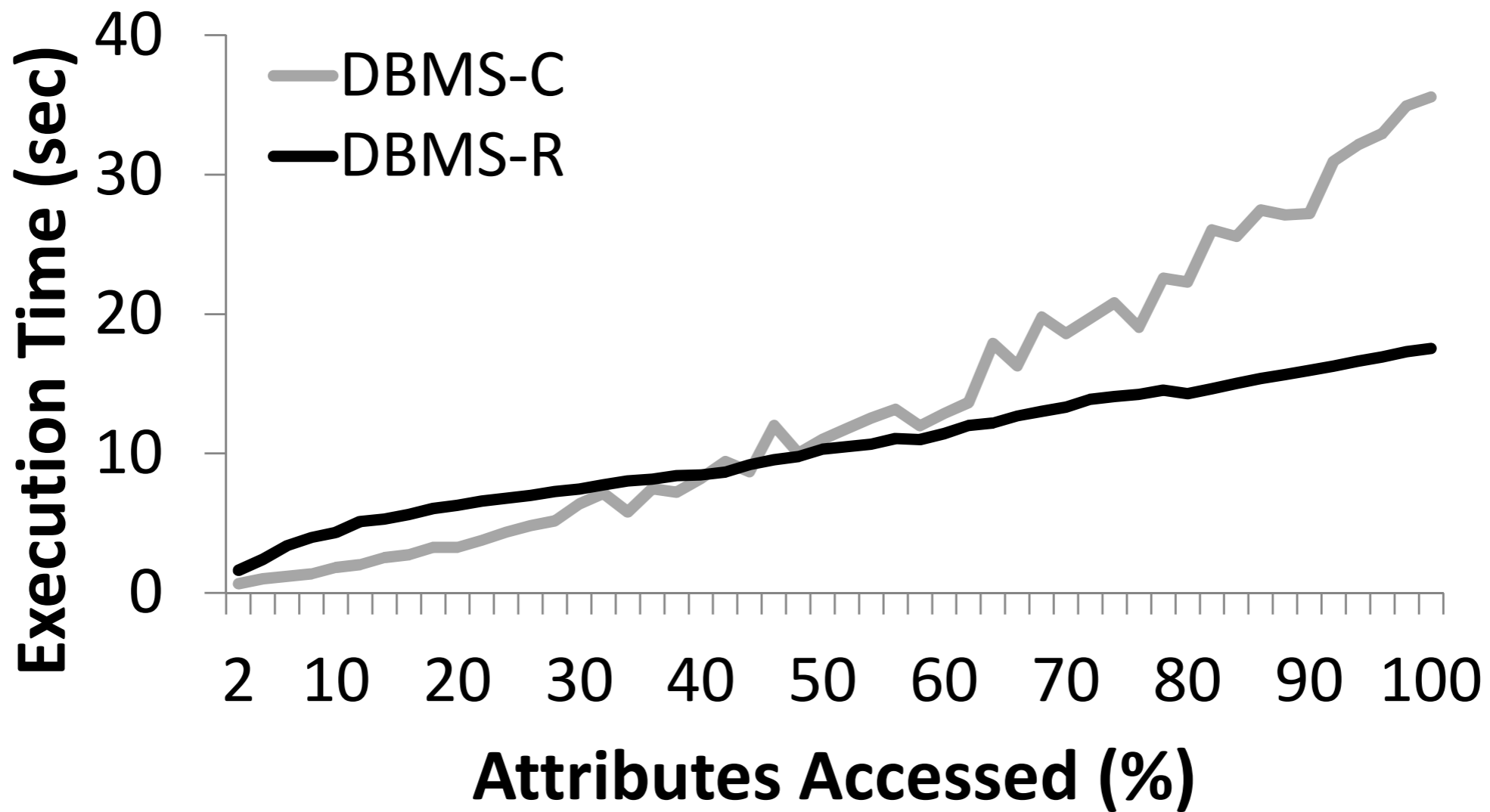


column-store

A B C D



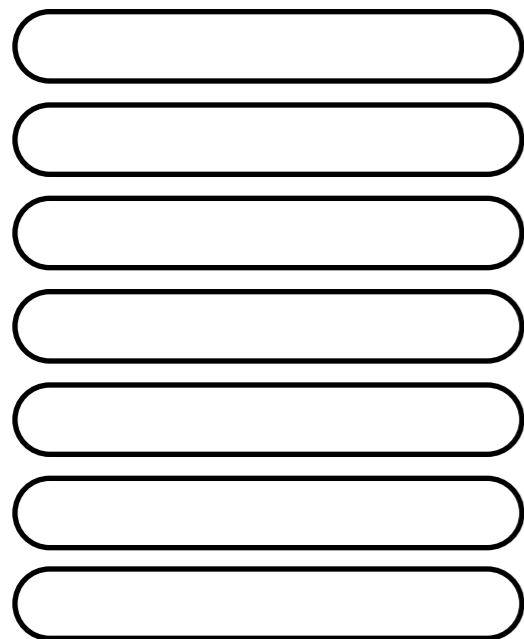
no fixed optimal solution



rows & columns

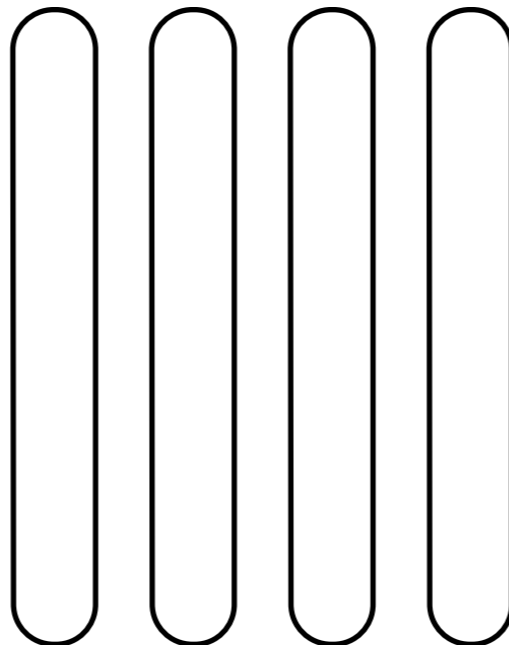
row-store

A B C D



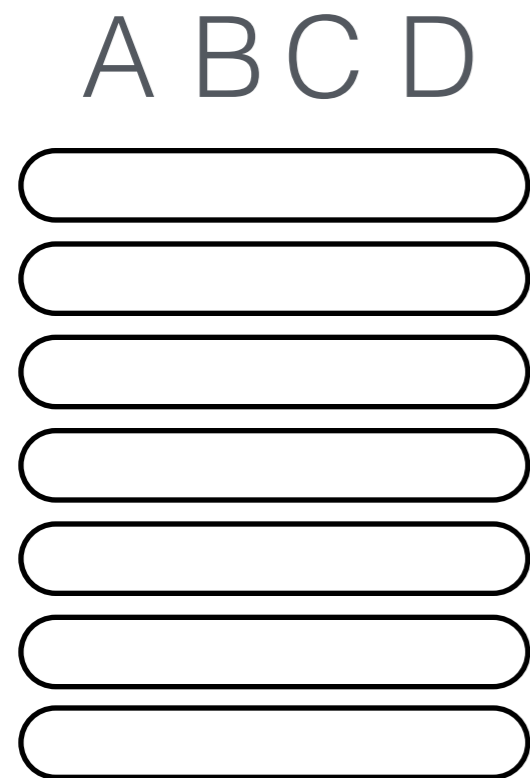
column-store

A B C D

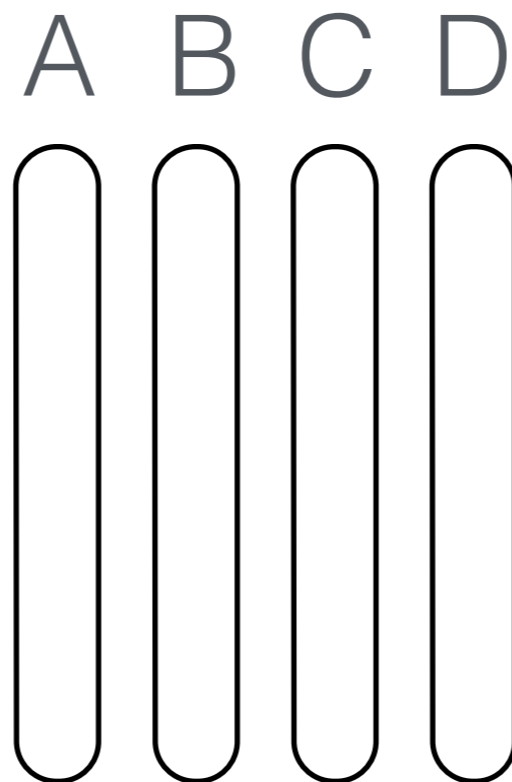


rows & columns

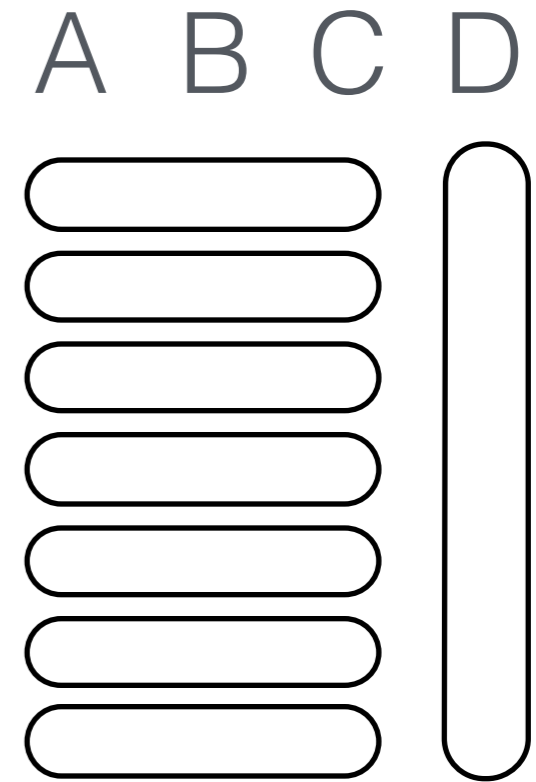
row-store



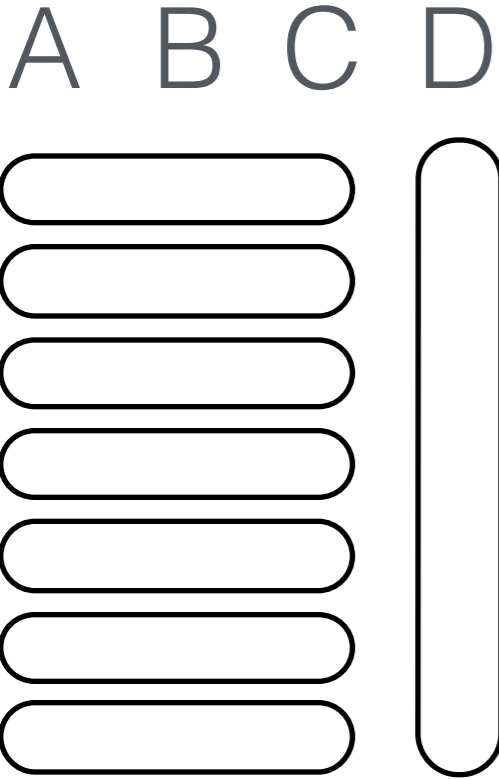
column-store



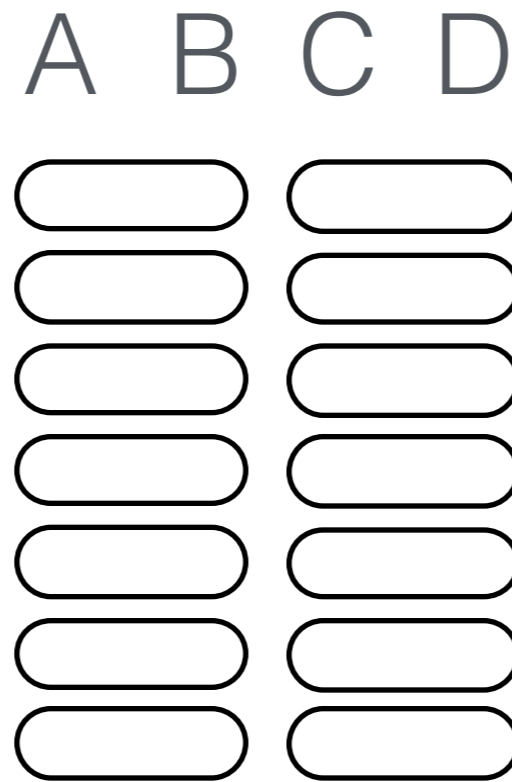
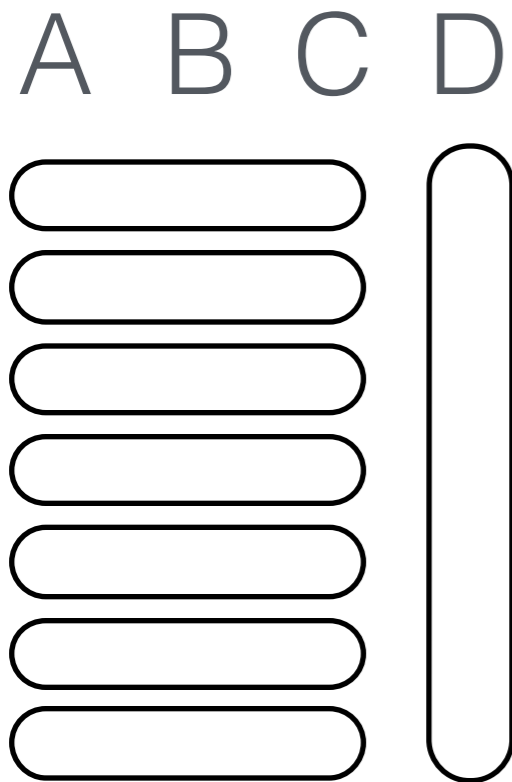
hybrid-store



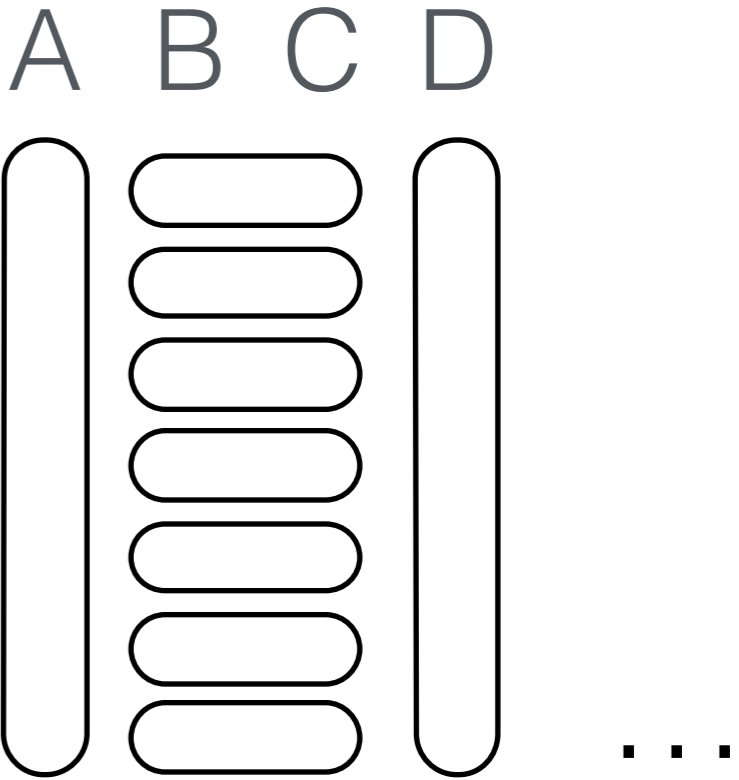
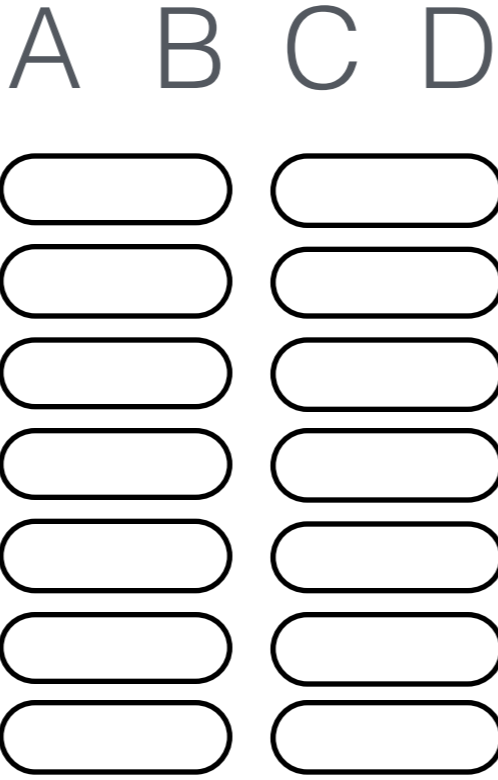
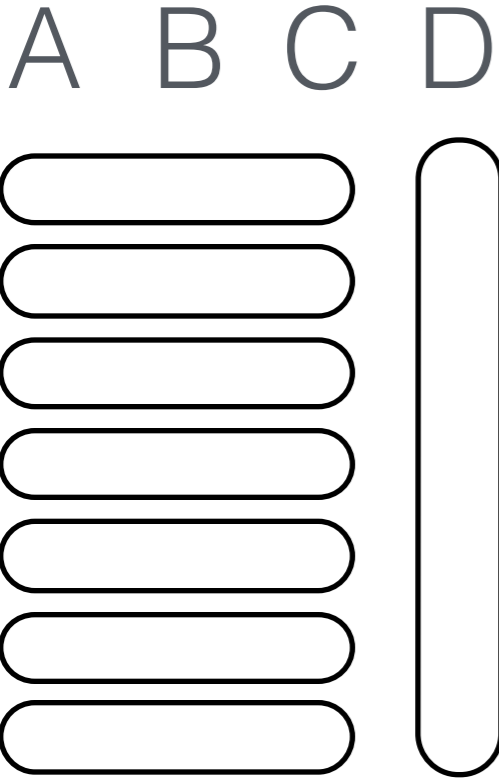
which layout is best?



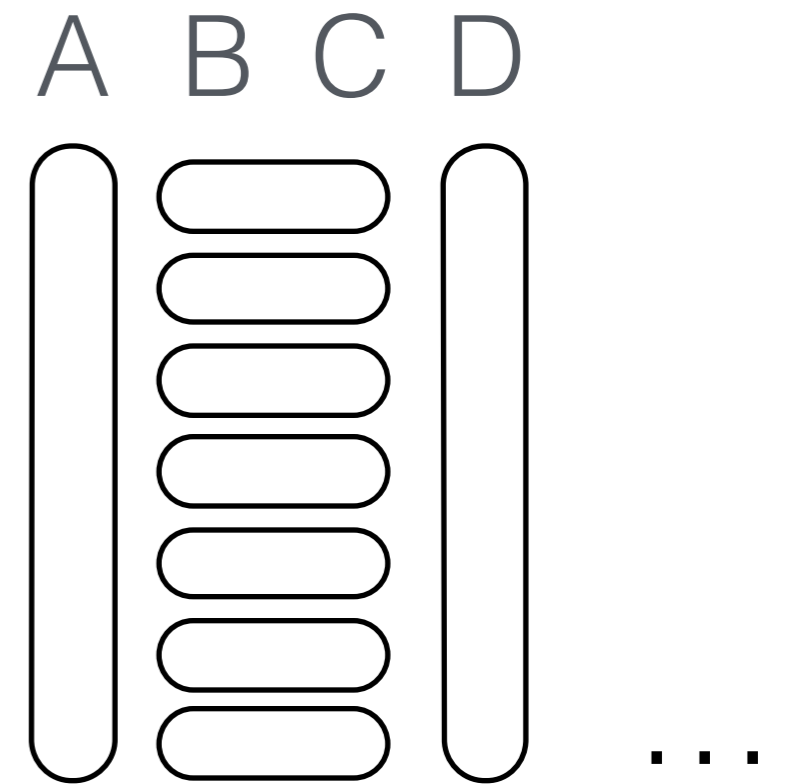
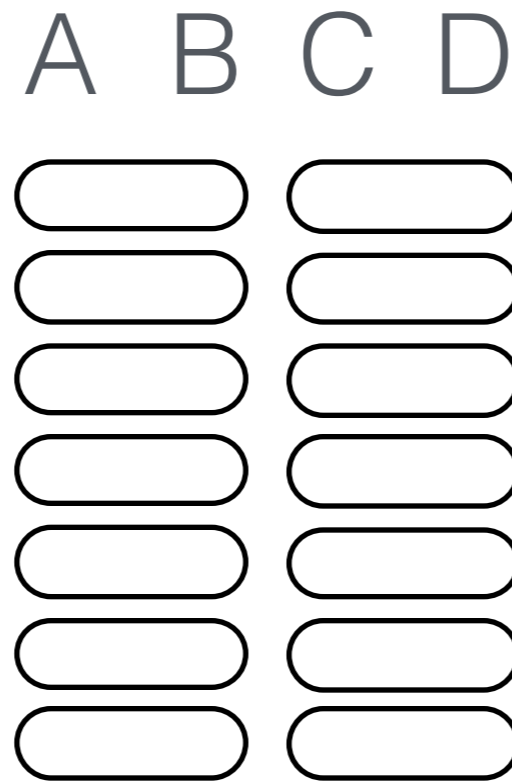
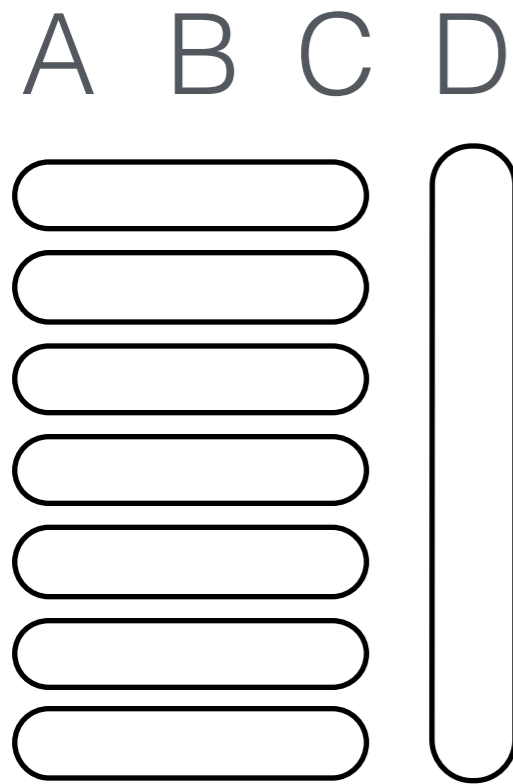
which layout is best?



which layout is best?



which layout is best?



too many combinations to maintain in parallel

query cost

$$q(L) = \sum_{i=1}^{|L|} \max(cost_i^{IO}, cost_i^{CPU})$$

for a given query we can know exactly which layout is best
the one that will cause the fewer cache misses

if we know all queries up front we can choose the layouts

adaptive storage:

continuously adapt layouts based on incoming queries

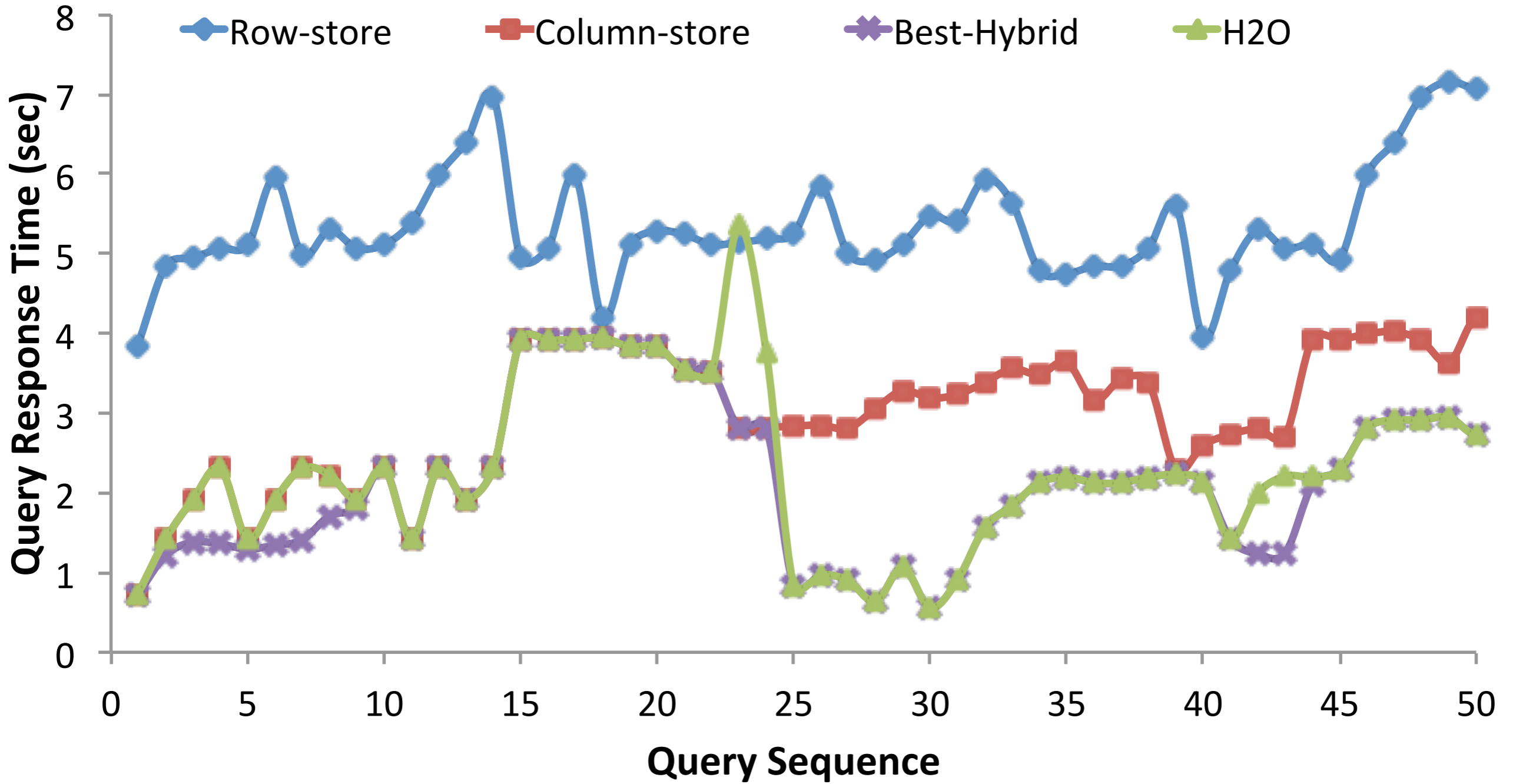
but computing all possible combinations is expensive...

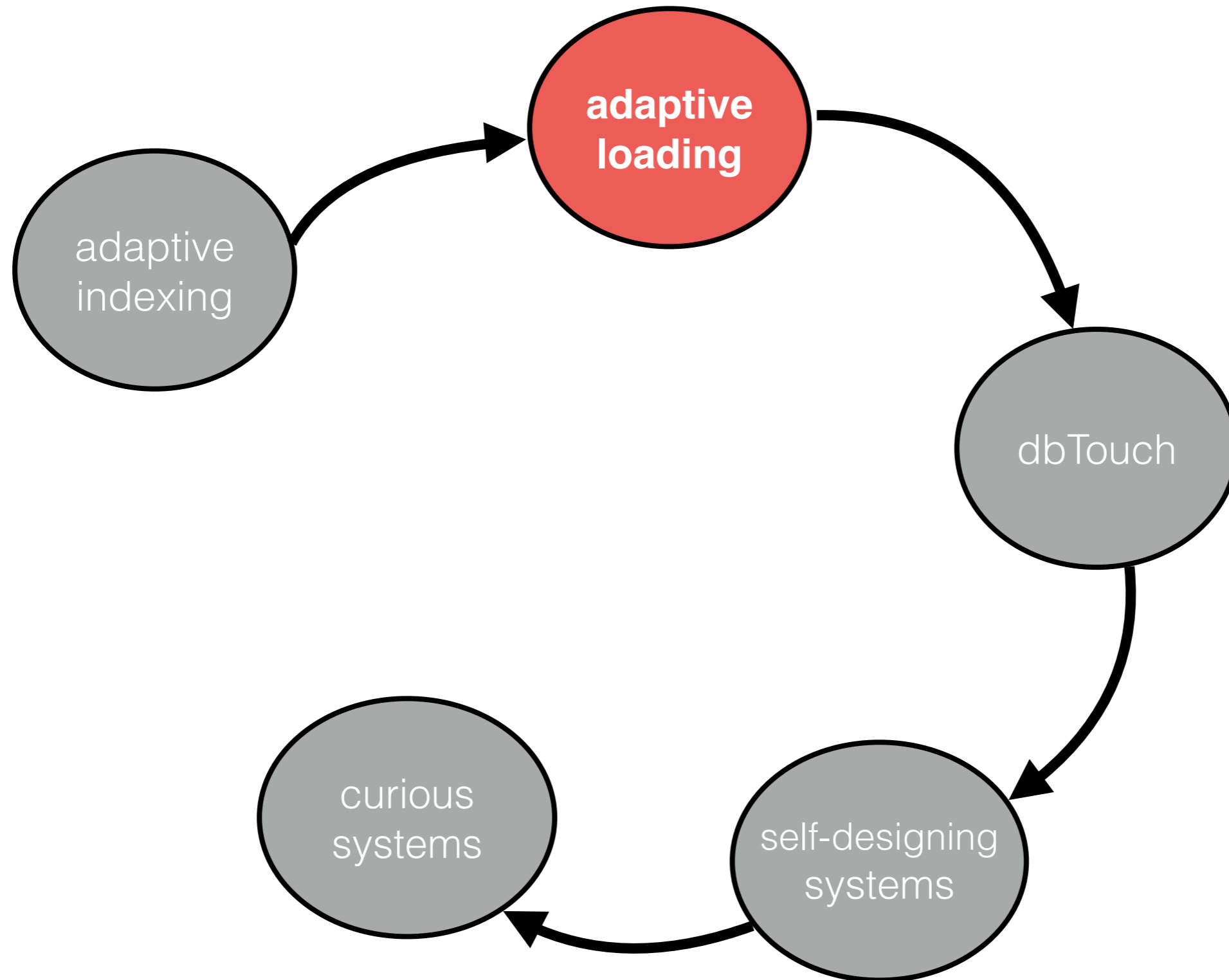


query

select A+B+C+D from R where A<10 and E>10

1. deal only with attributes referenced in queries
2. handle select clause separately from where clause
3. start from pure column-store and build up
4. stop when no improvement possible





loading



loading



copy data inside the database
database now has full control

loading

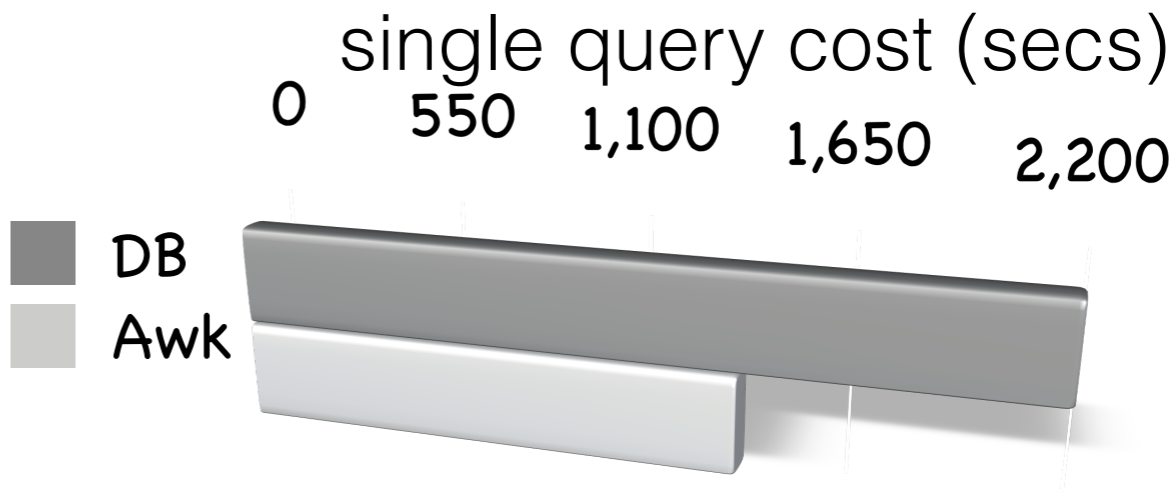


copy data inside the database
database now has full control

slow process...
not all data might be needed all the time

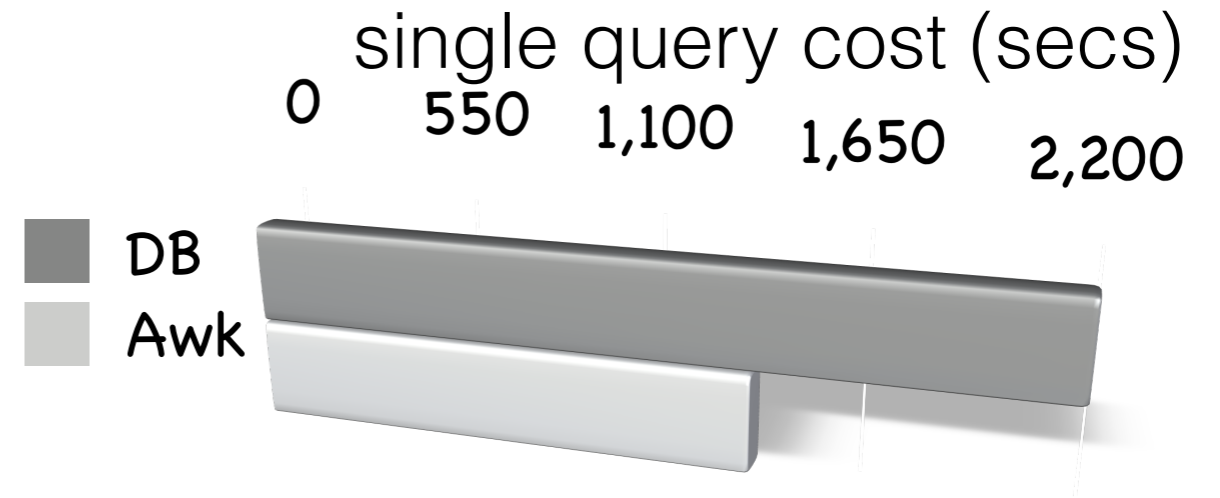
database vs. unix tools

1 file, 4 attributes,
1 billion tuples



database vs. unix tools

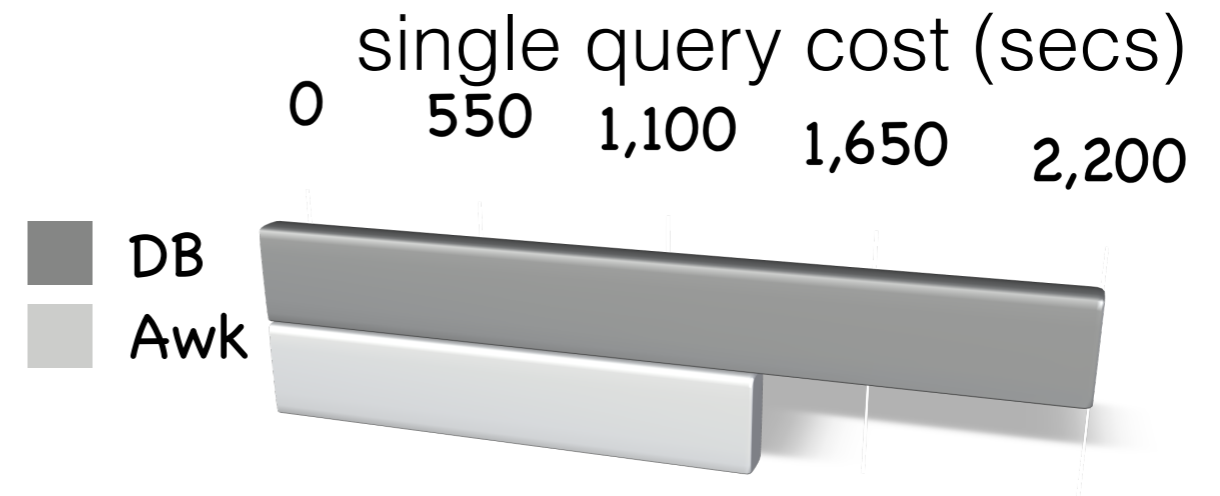
1 file, 4 attributes,
1 billion tuples



break down db cost

database vs. unix tools

1 file, 4 attributes,
1 billion tuples

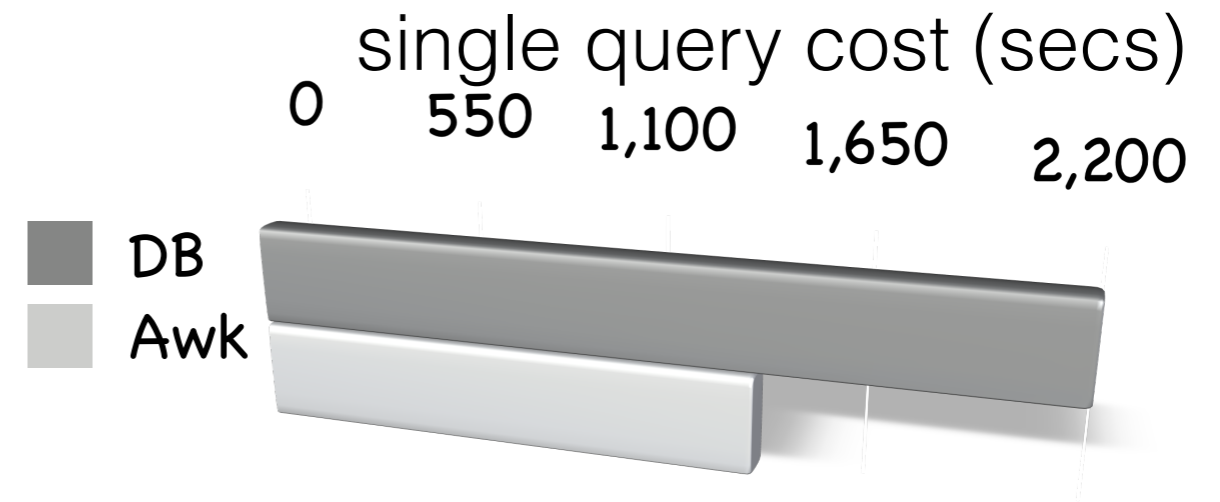


break down db cost

loading is a major bottleneck

database vs. unix tools

1 file, 4 attributes,
1 billion tuples



break down db cost

loading is a major bottleneck

but writing/maintaining scripts does not scale

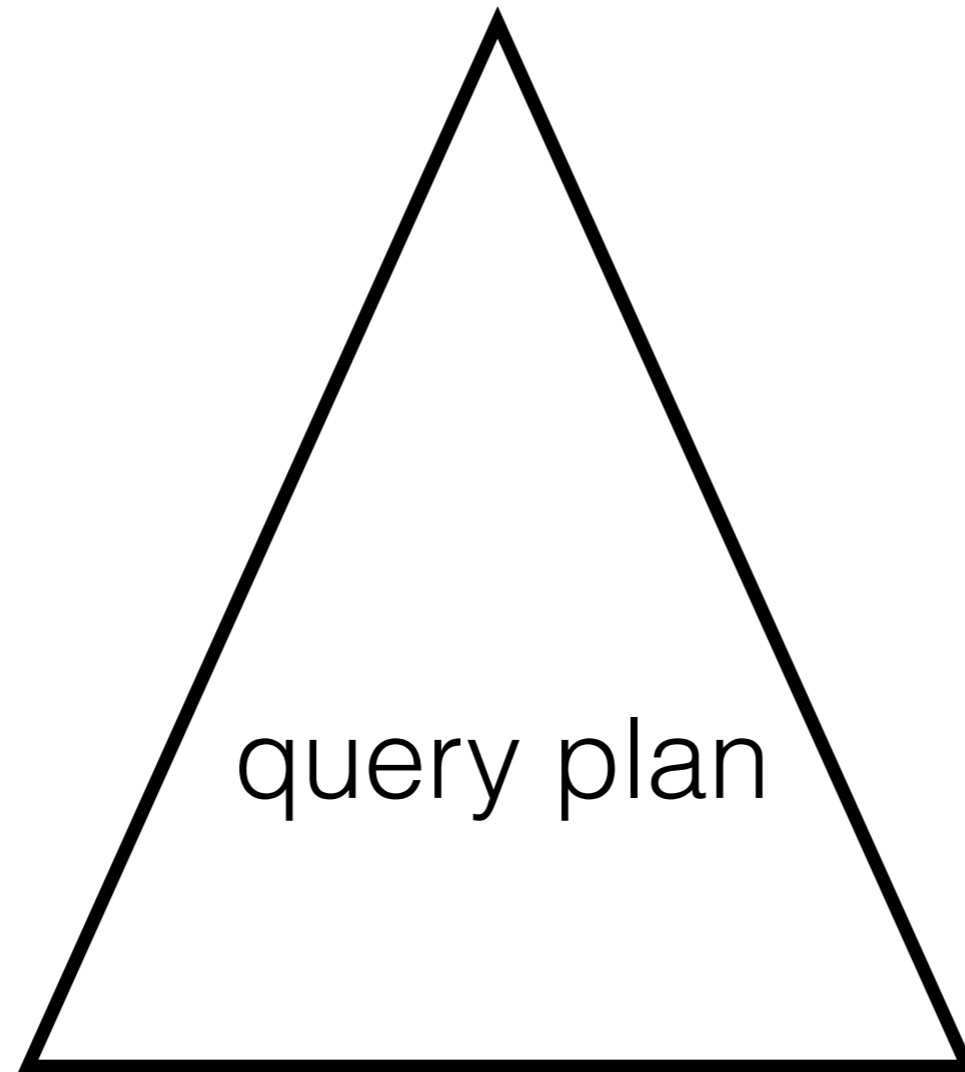
adaptive loading

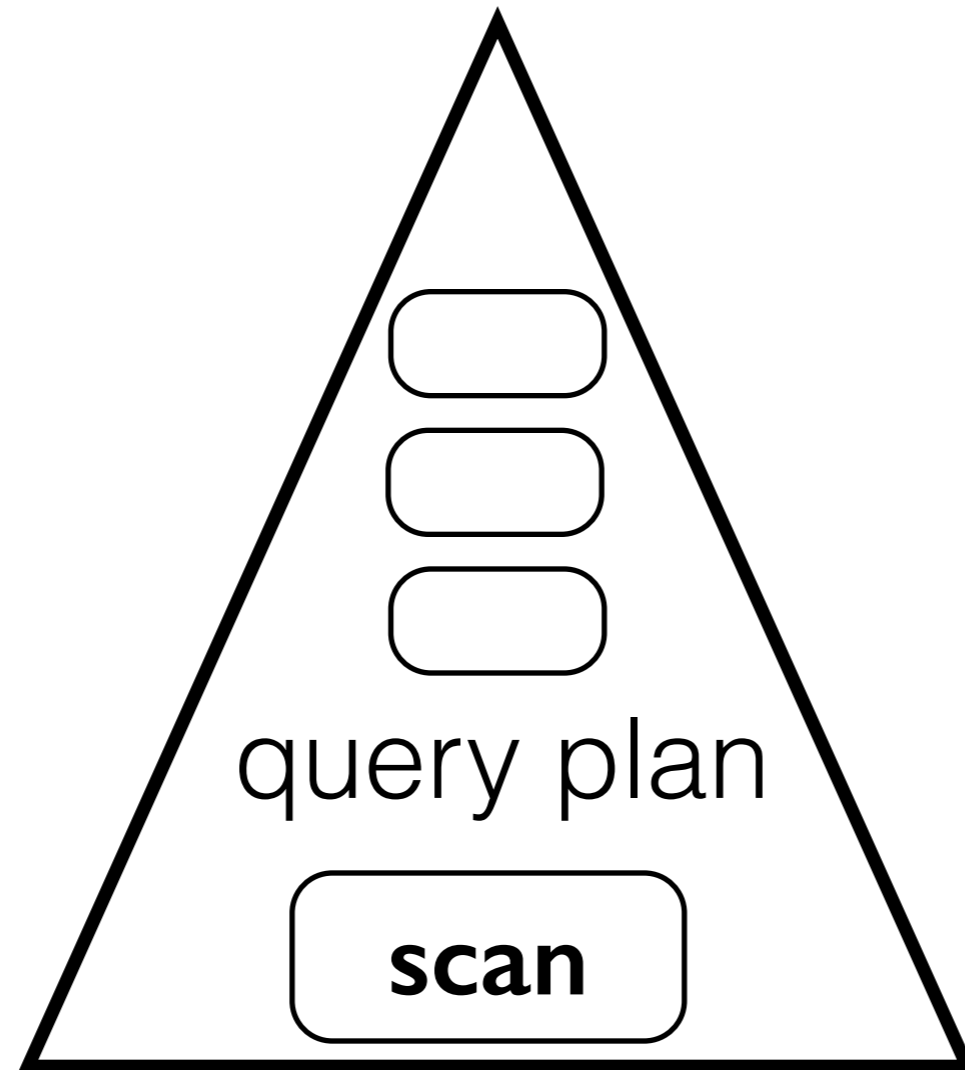
**load/touch only what is needed
and only when it is needed**

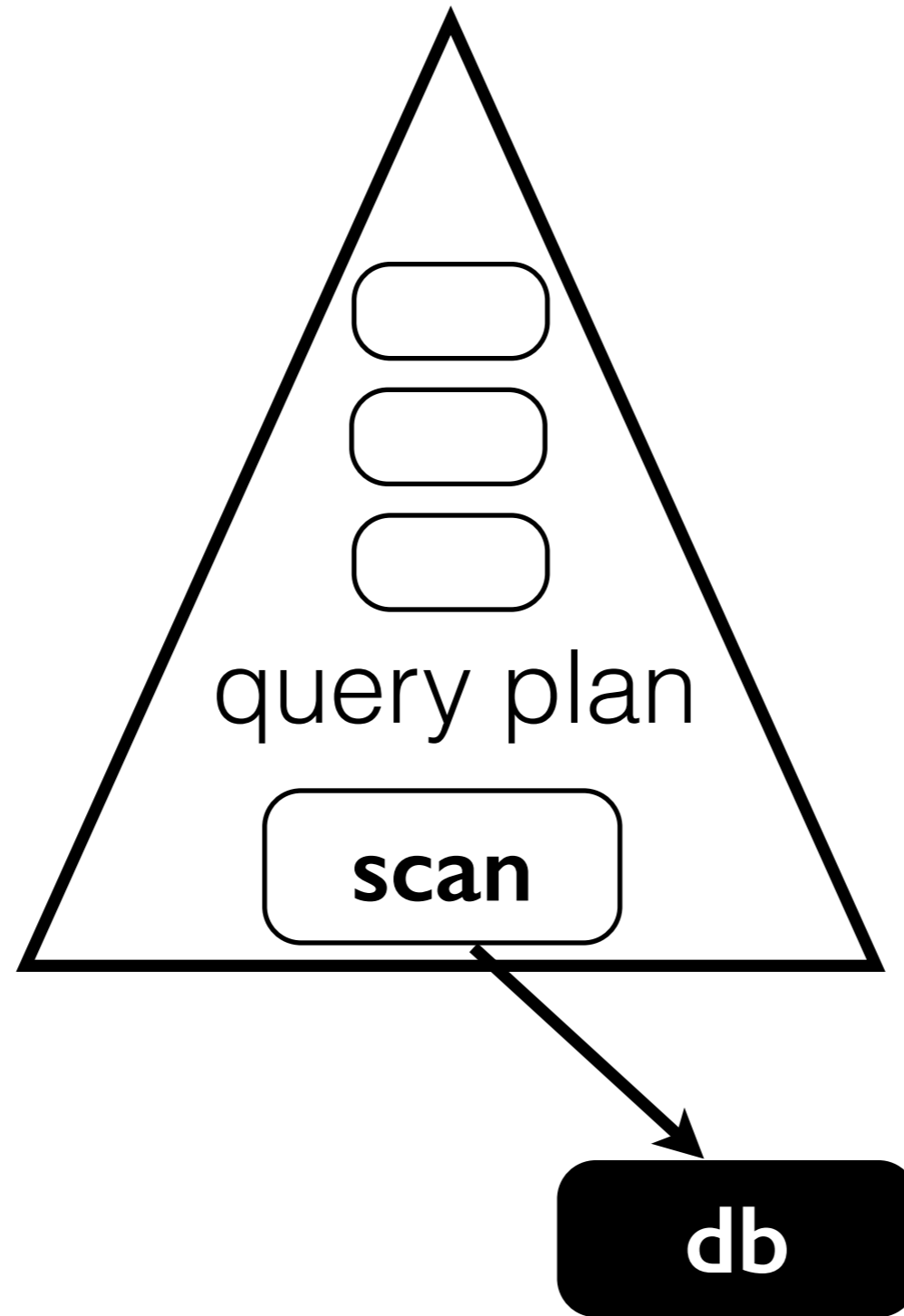
but raw data access is expensive

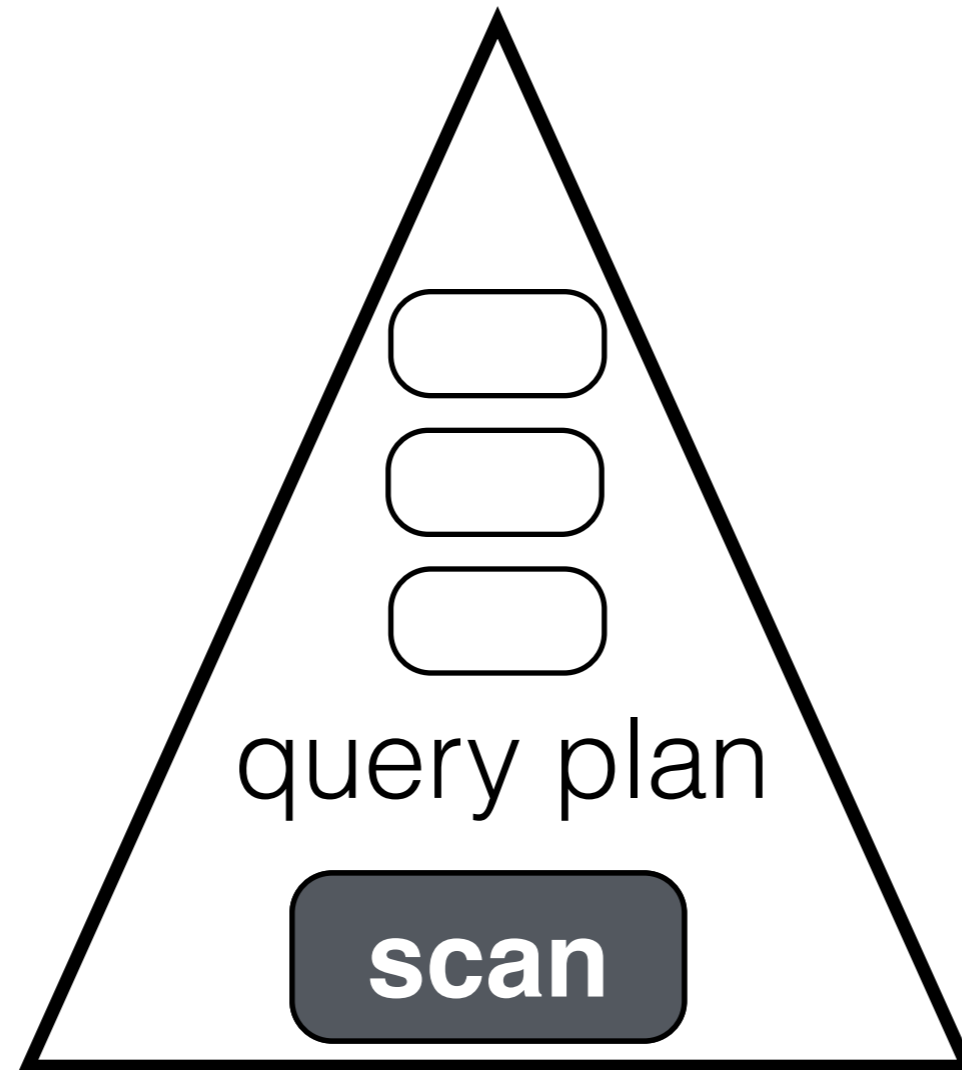
tokenizing - parsing - no indexing - no statistics

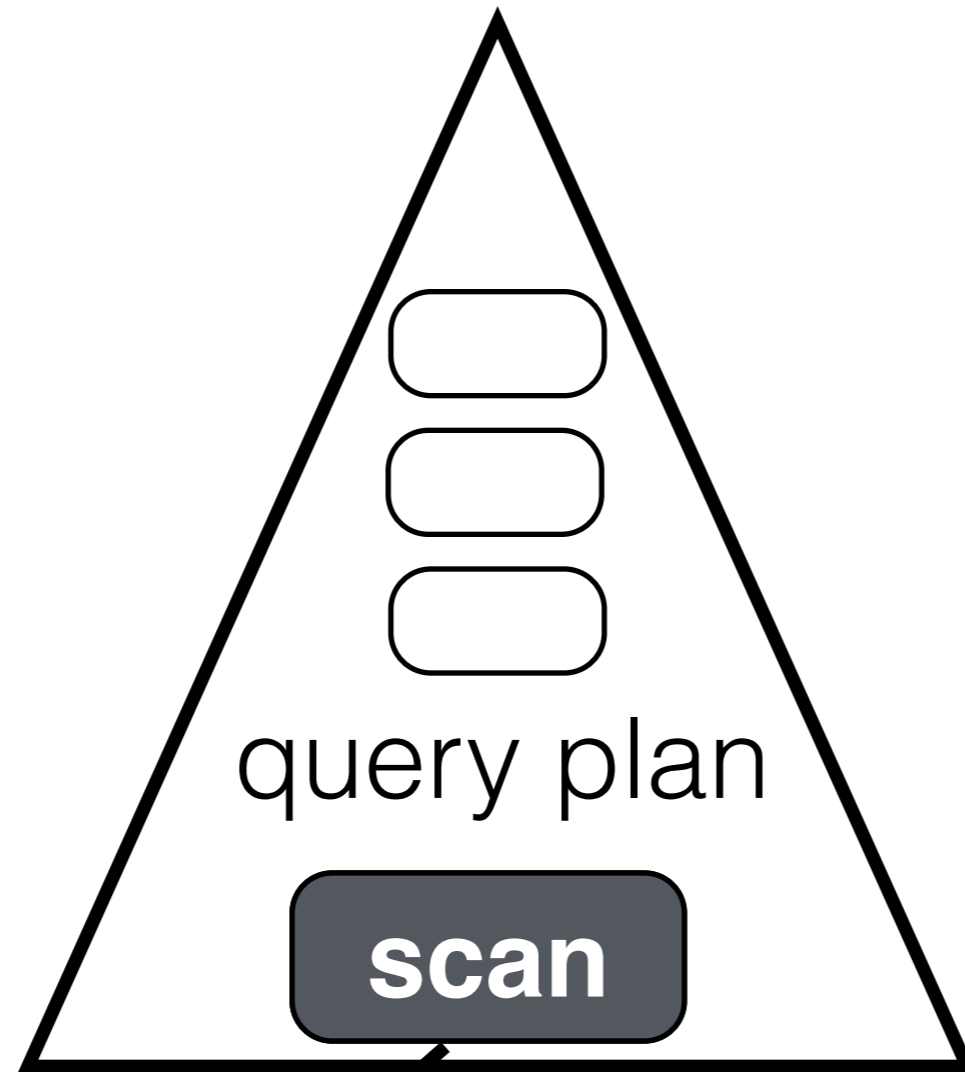
challenge: fast raw data access







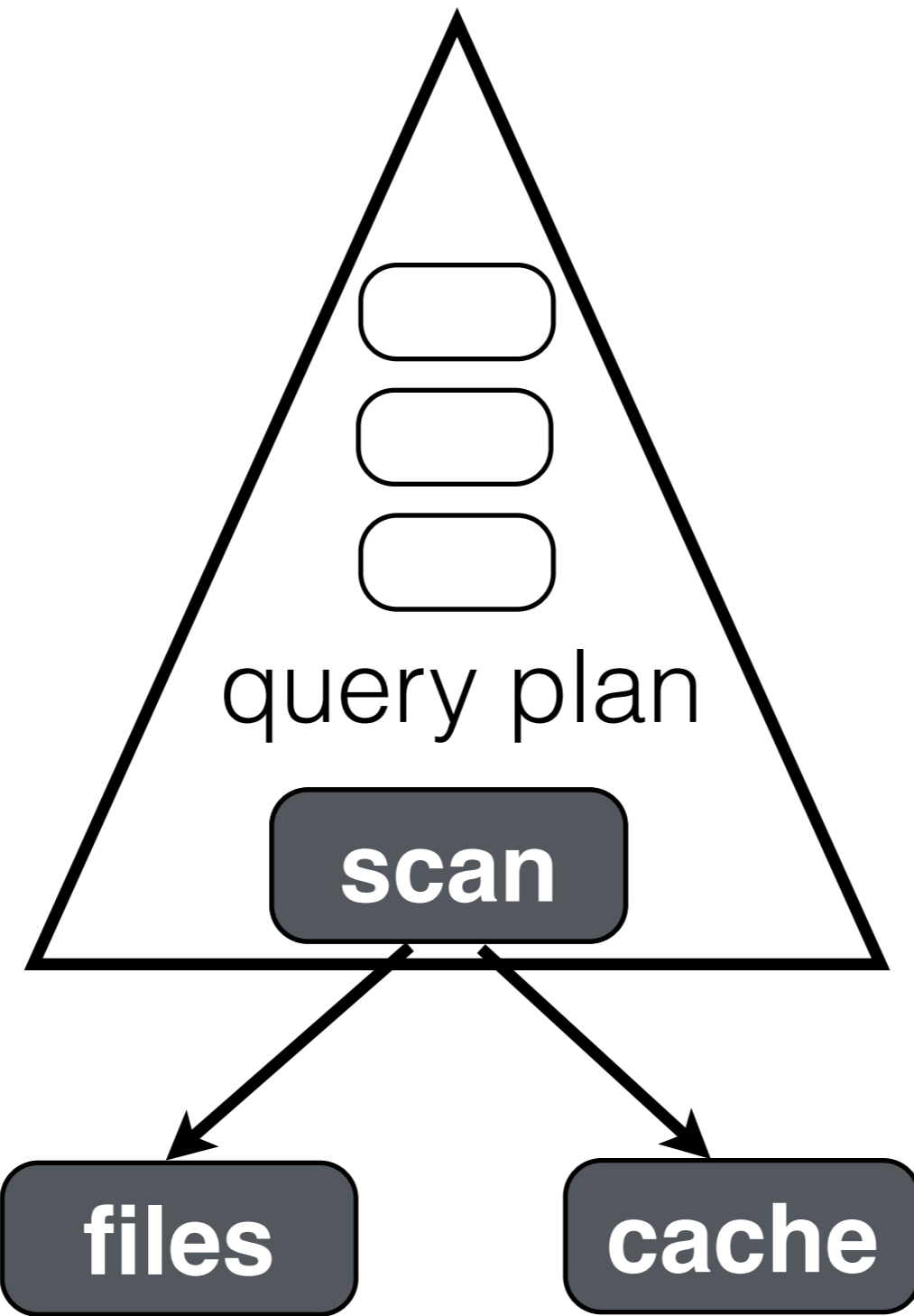




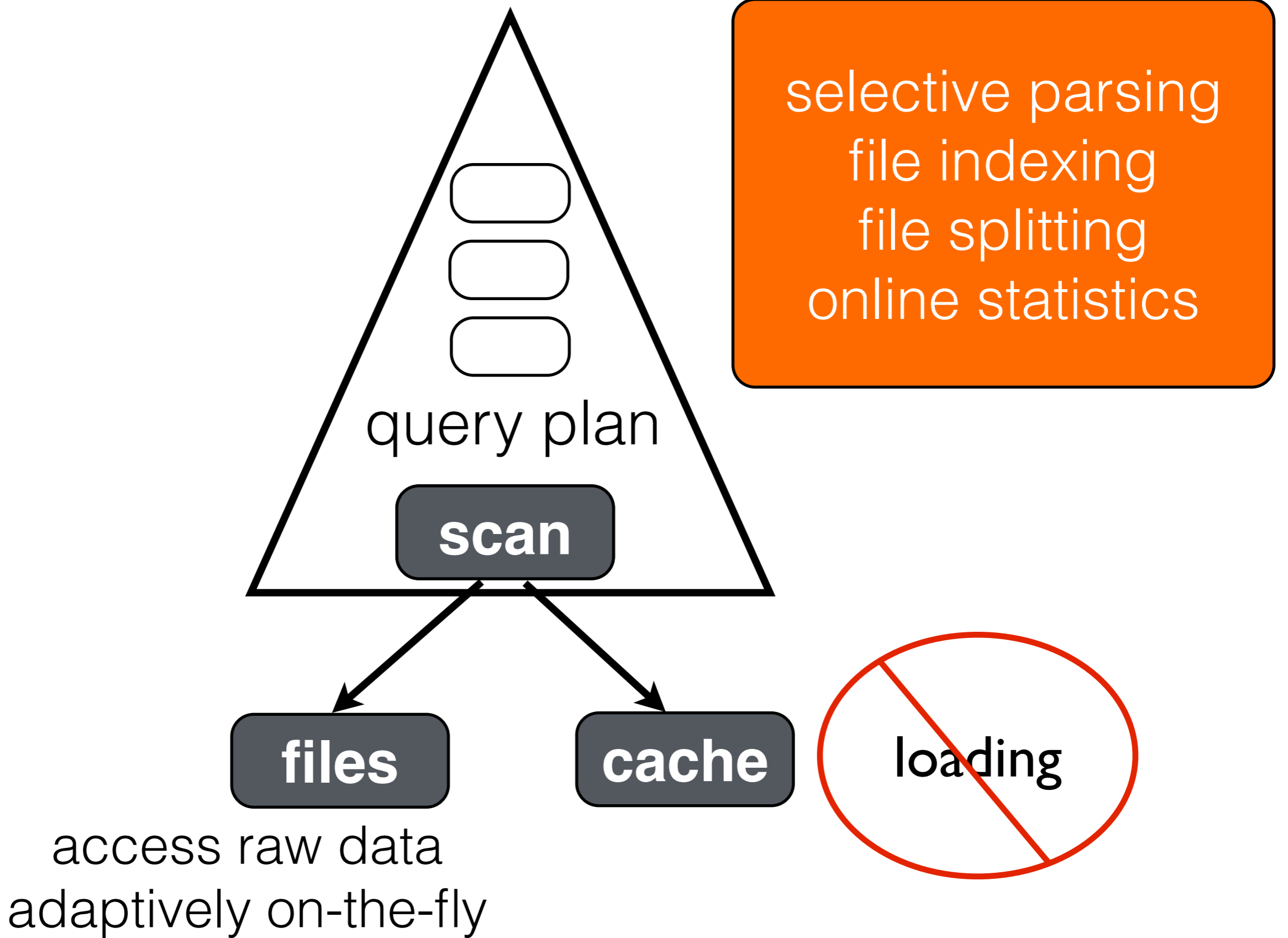
files

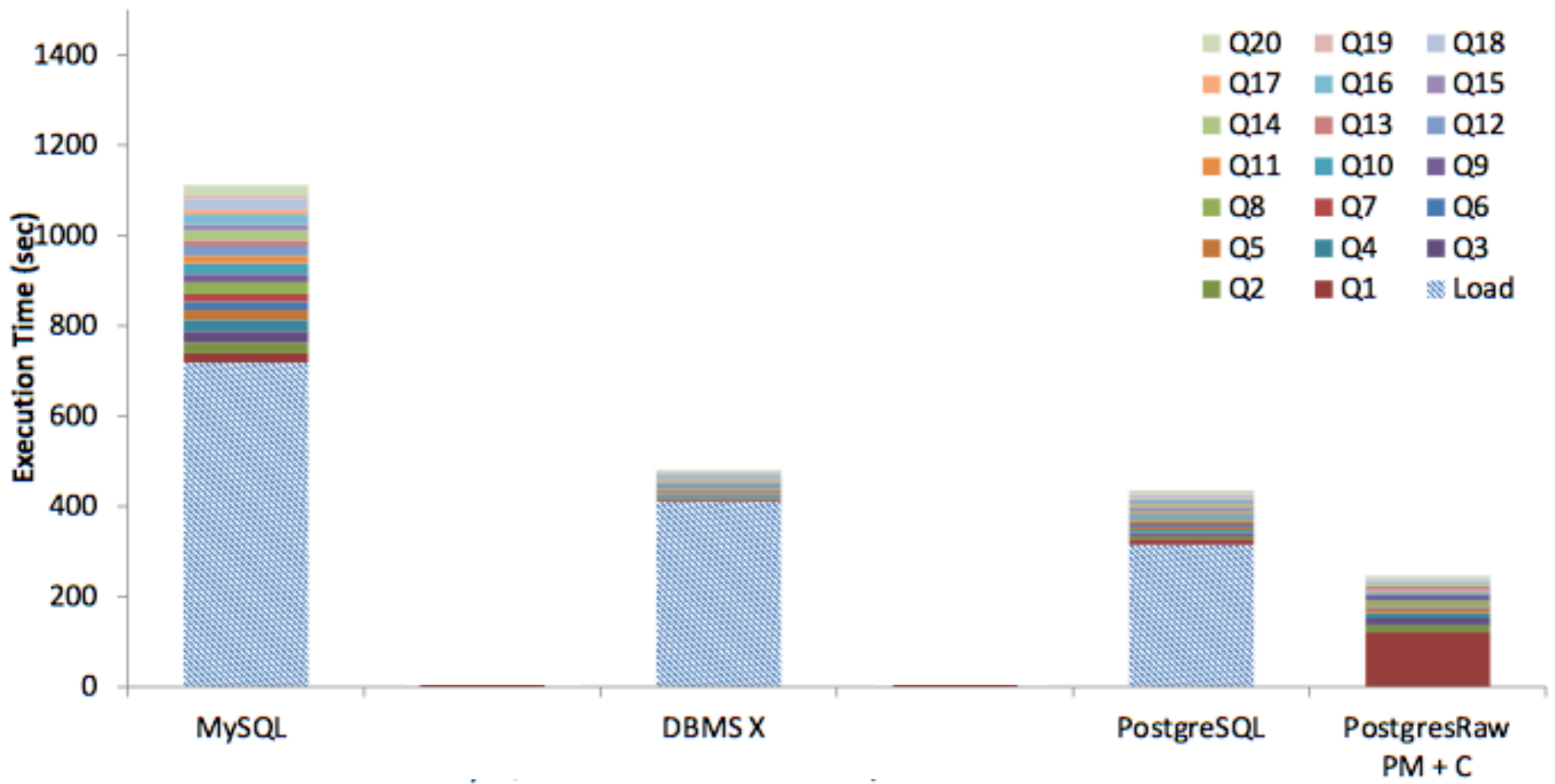
access raw data
adaptively on-the-fly

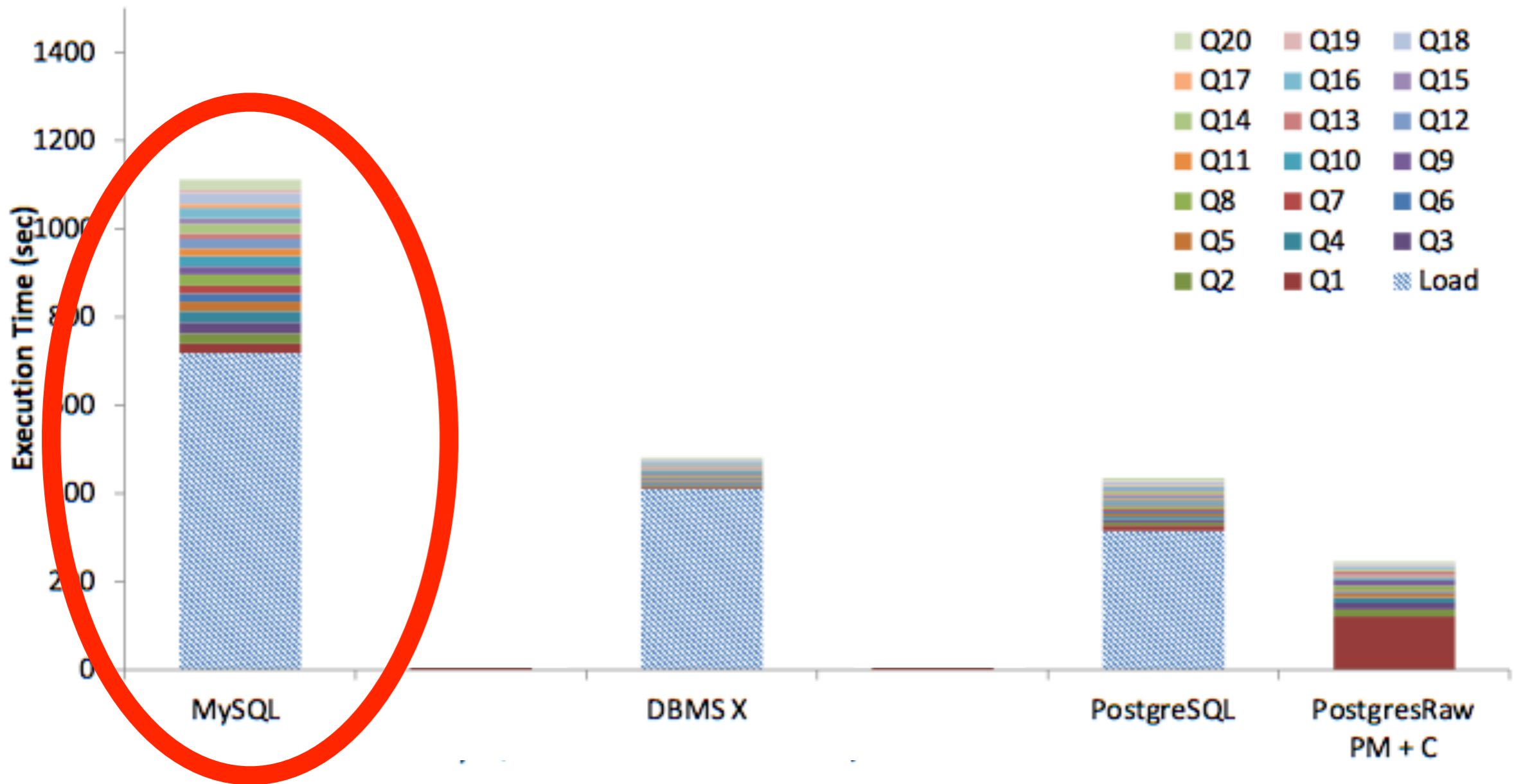
~~loading~~

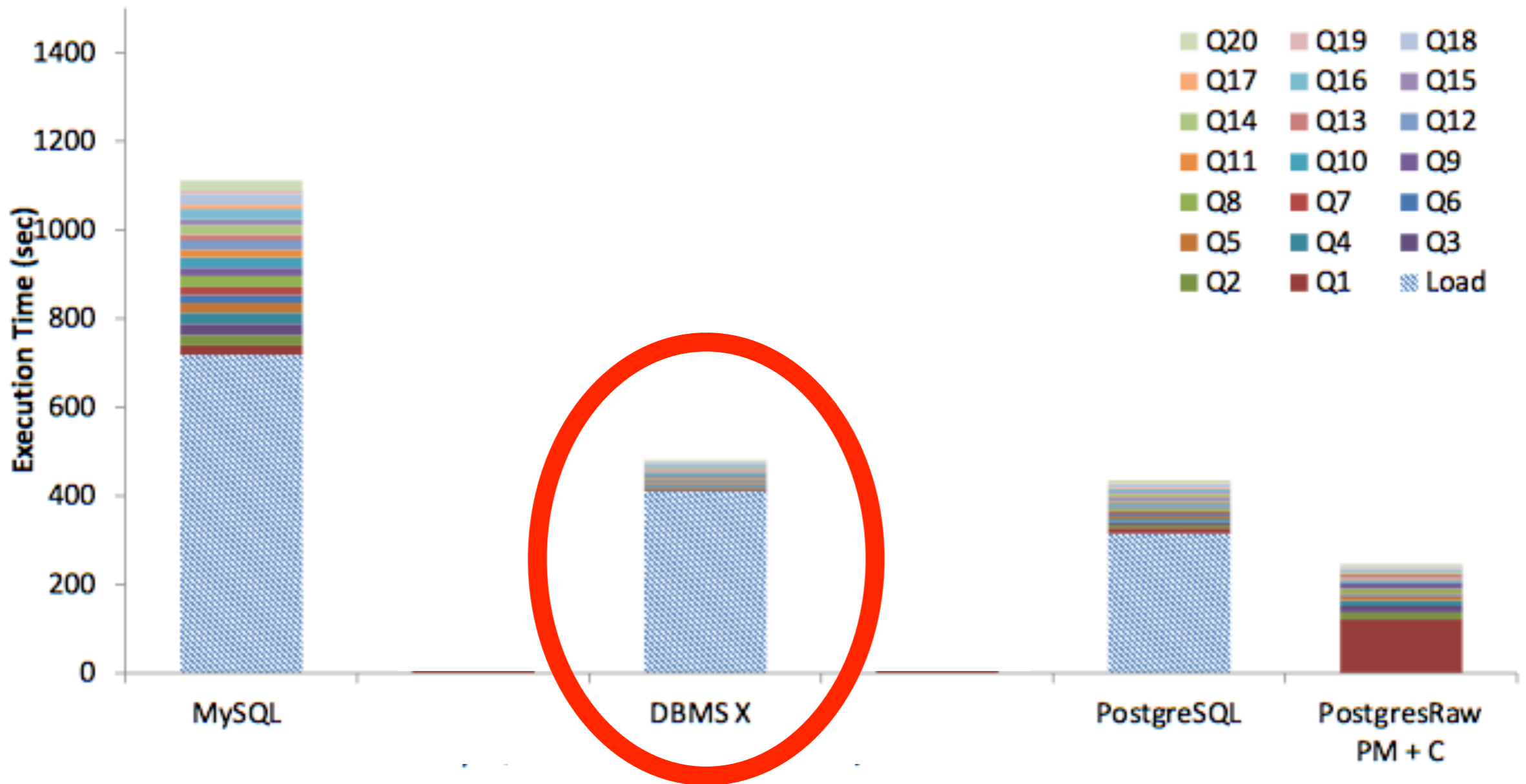


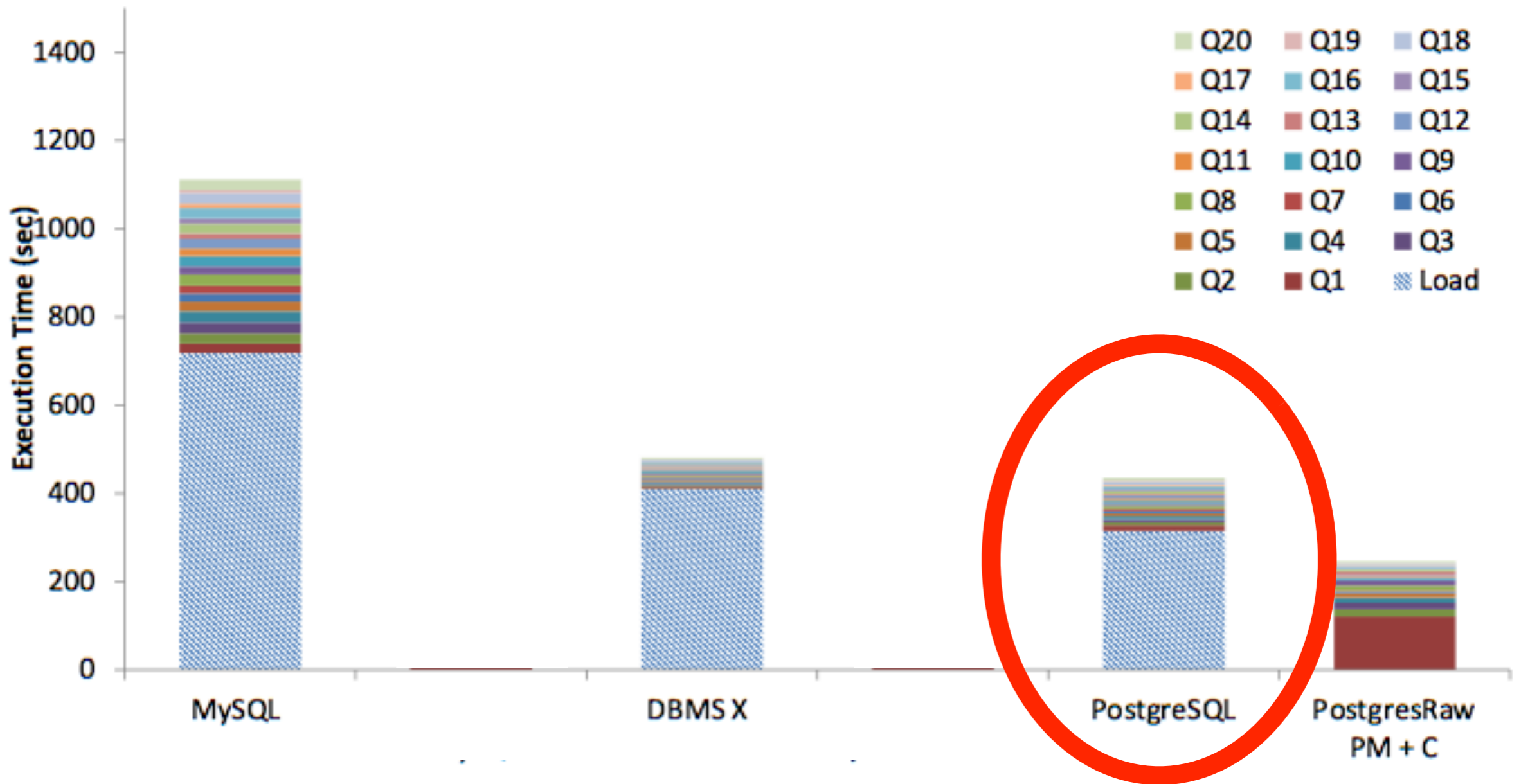
access raw data
adaptively on-the-fly

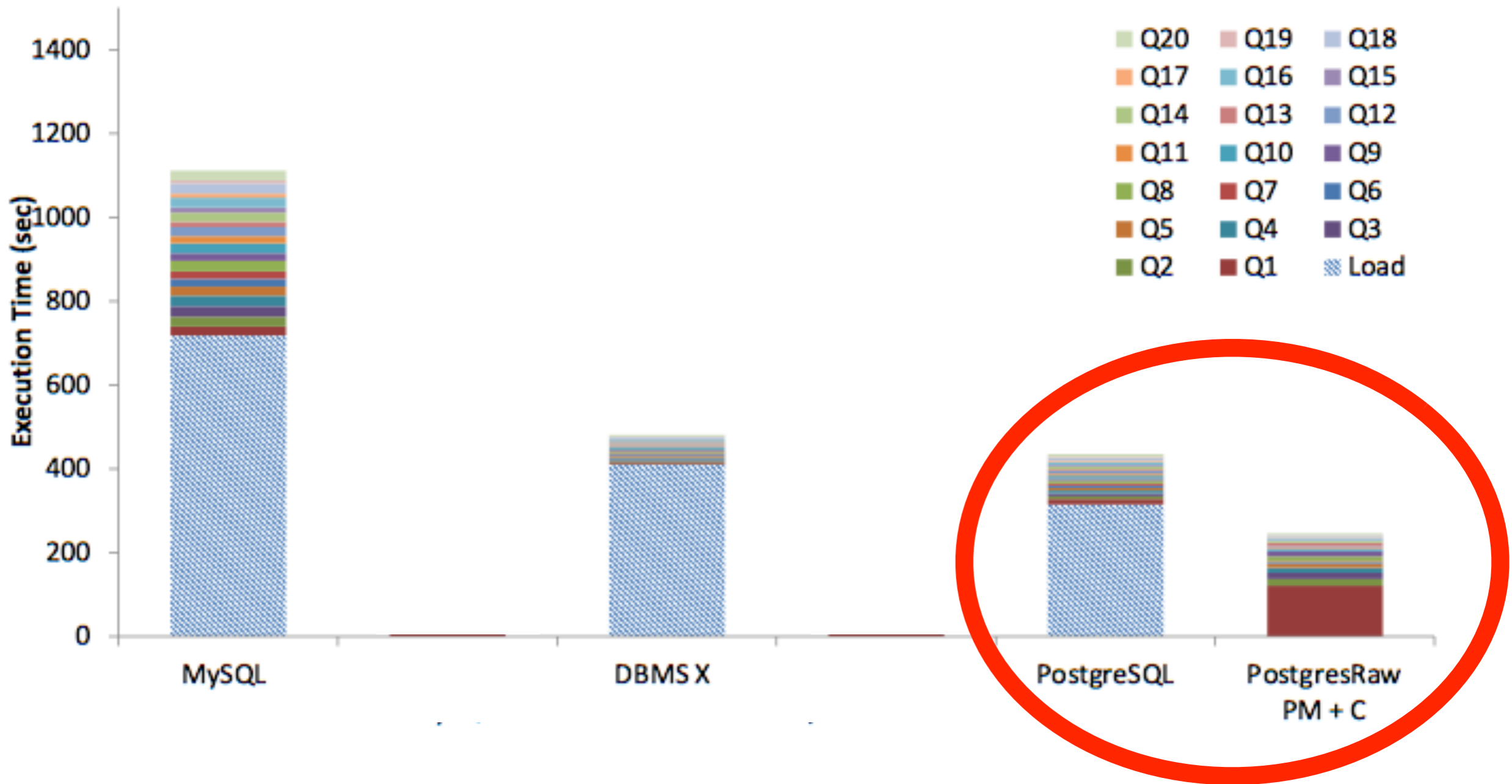




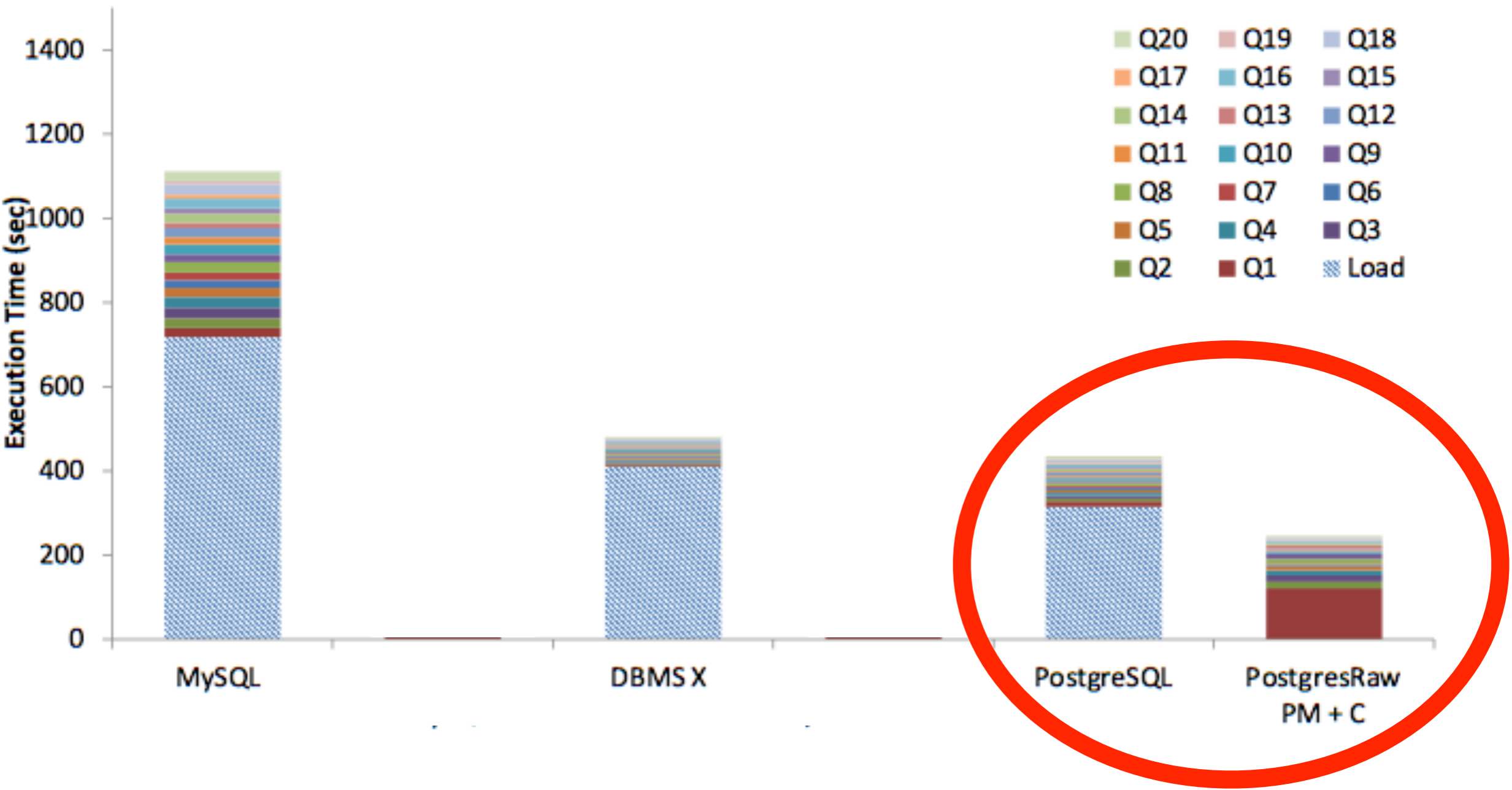


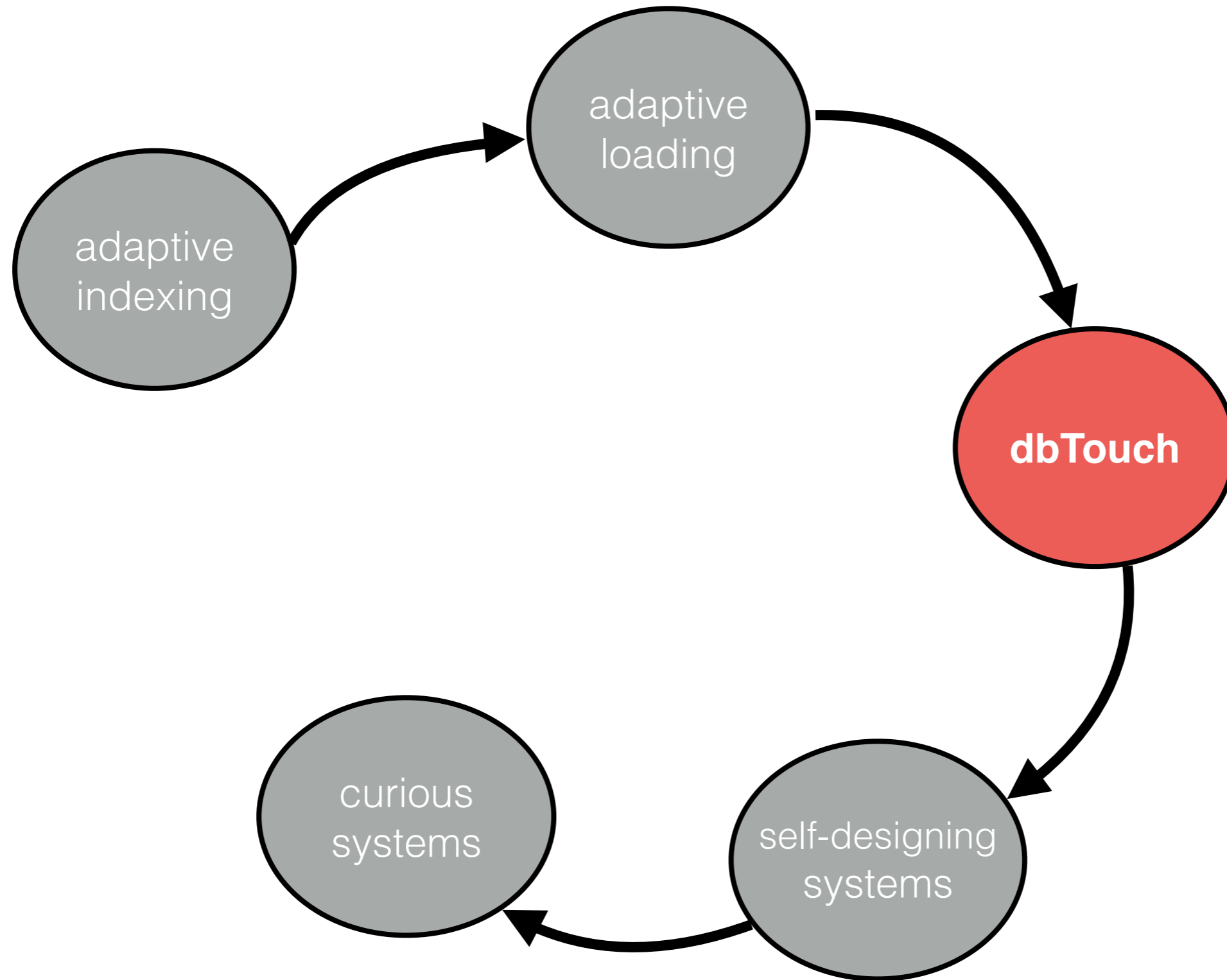






reducing data-to-query time





querying



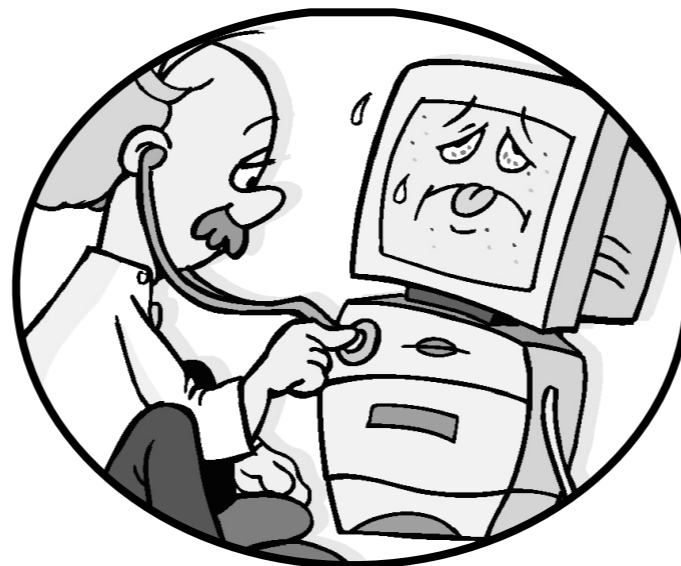
querying

load

tune

query

SQL interface



querying

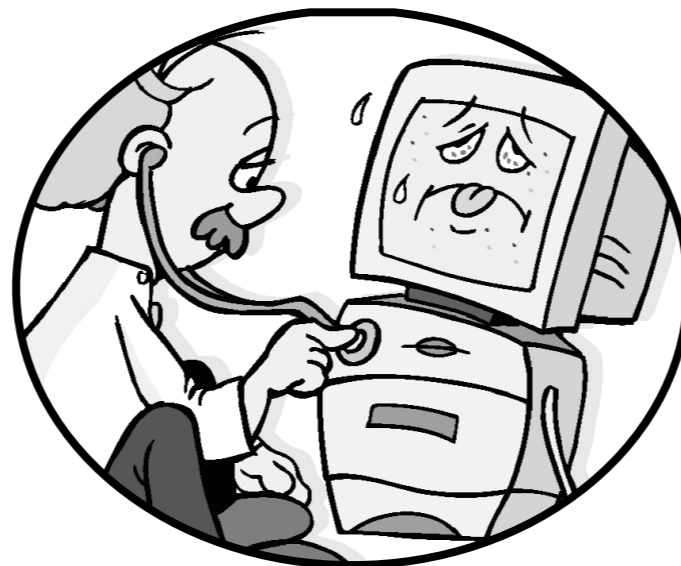
load

tune

query

SQL interface

correct and complete answers

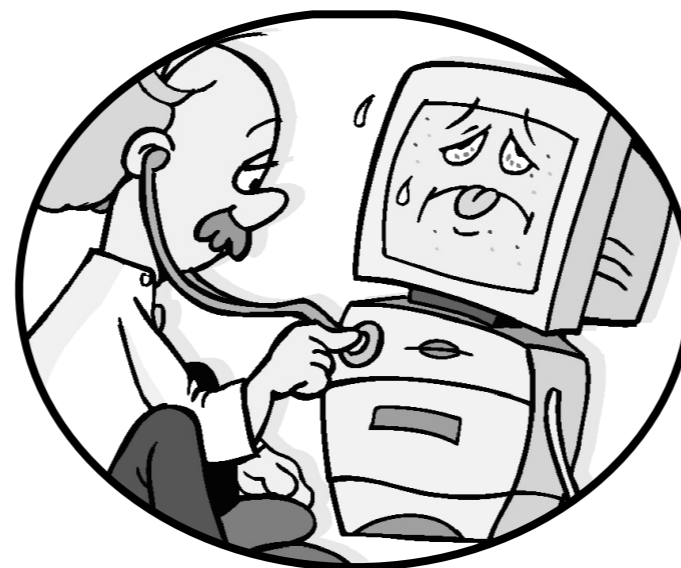
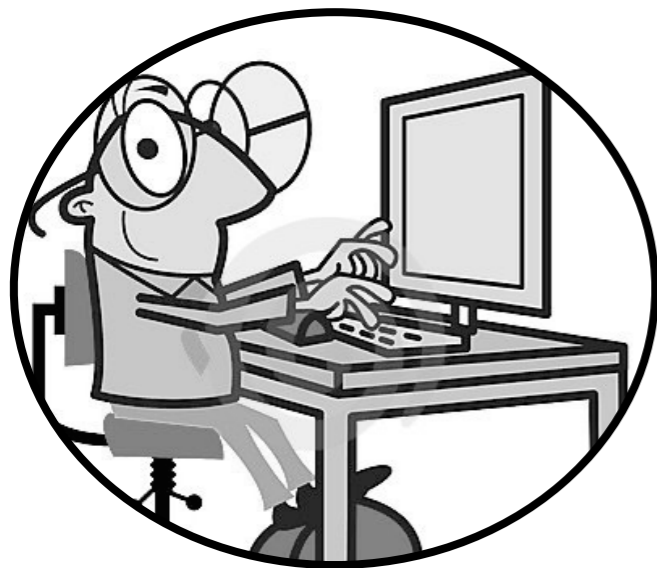


querying

complex and slow - not fit for exploration

SQL interface

correct and complete answers





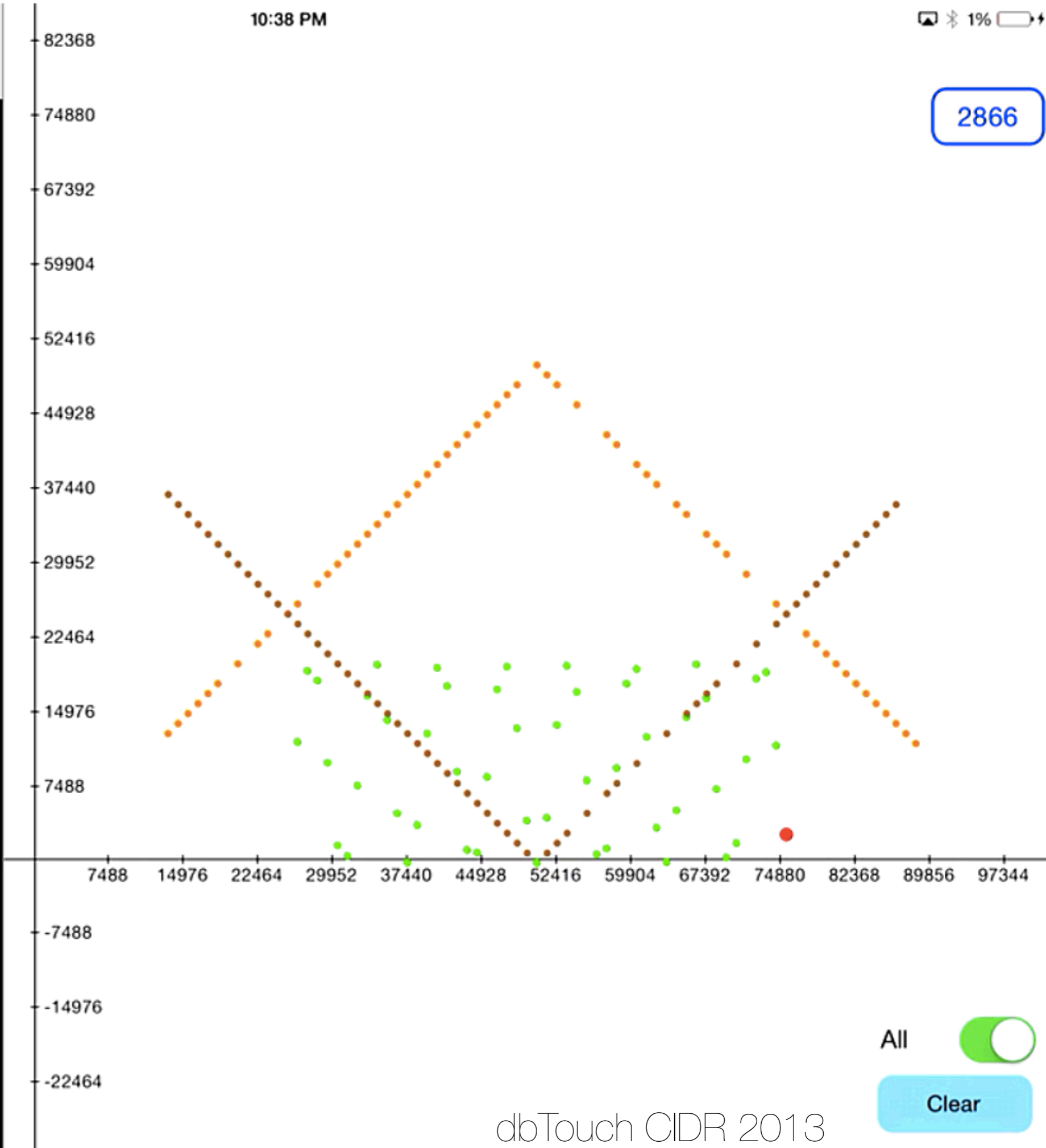
just touch the data you need



just touch the data you need

**this is not about query building
it is about query processing**

2866

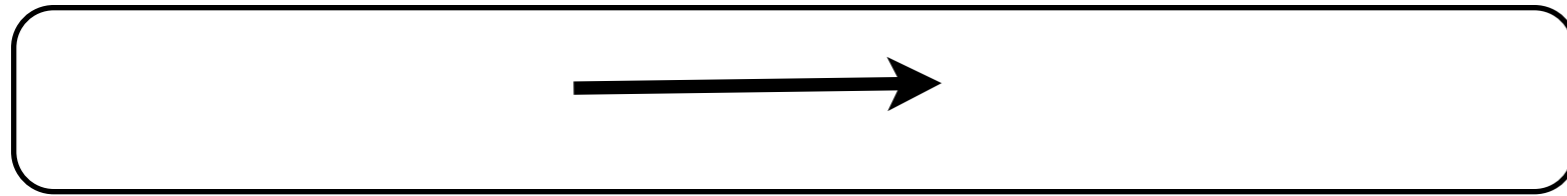


All

Clear

what does this mean for db kernels?

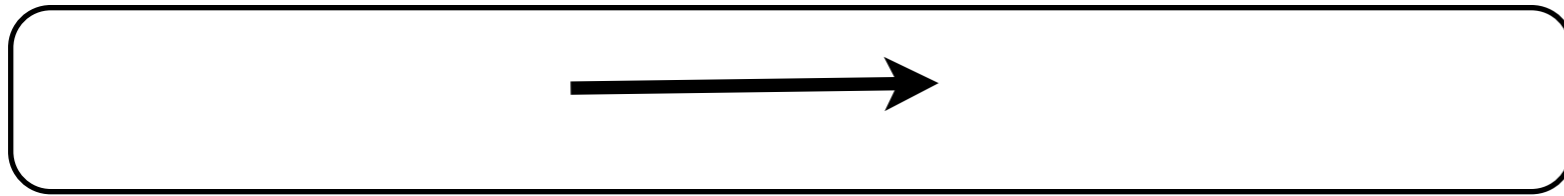
db



select R.a from R

what does this mean for db kernels?

db



select R.a from R

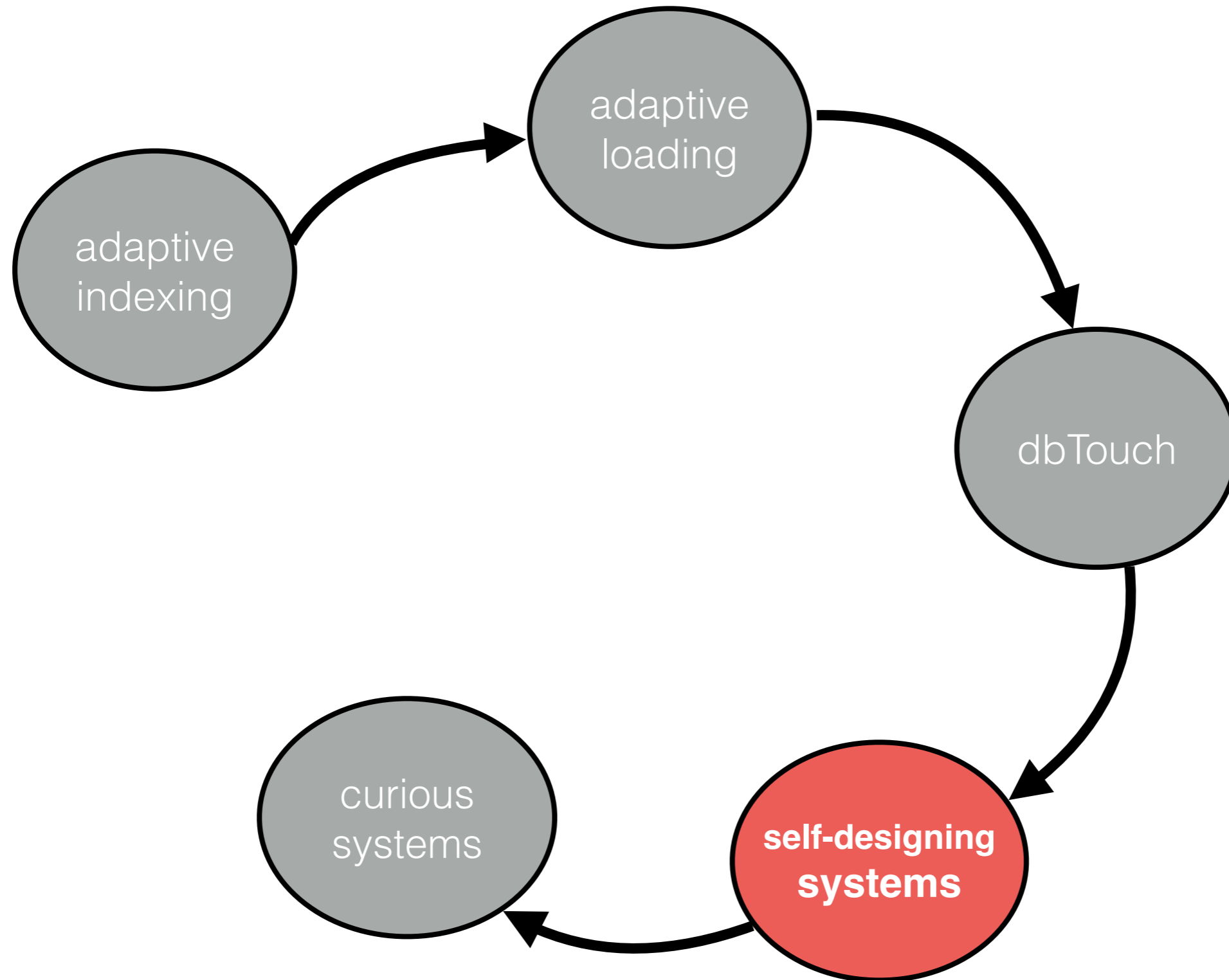
what does this mean for db kernels?

dbTouch

56 38 45 2

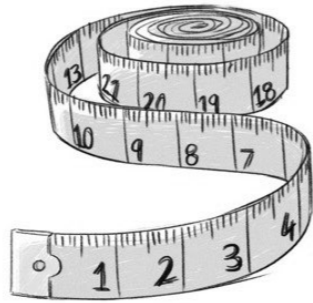


process only
what you touch





+



+



=



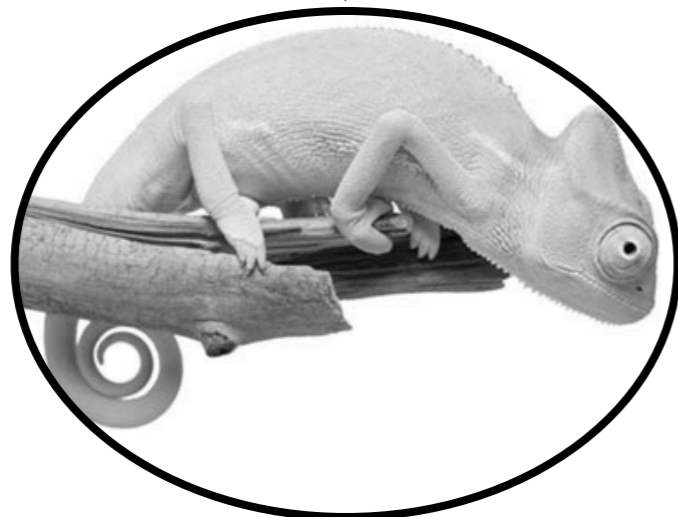
One size does **not** fit all

Custom solutions are needed for optimal performance

Solutions need to be tuned

Bootstrapping new systems is expensive and time-consuming

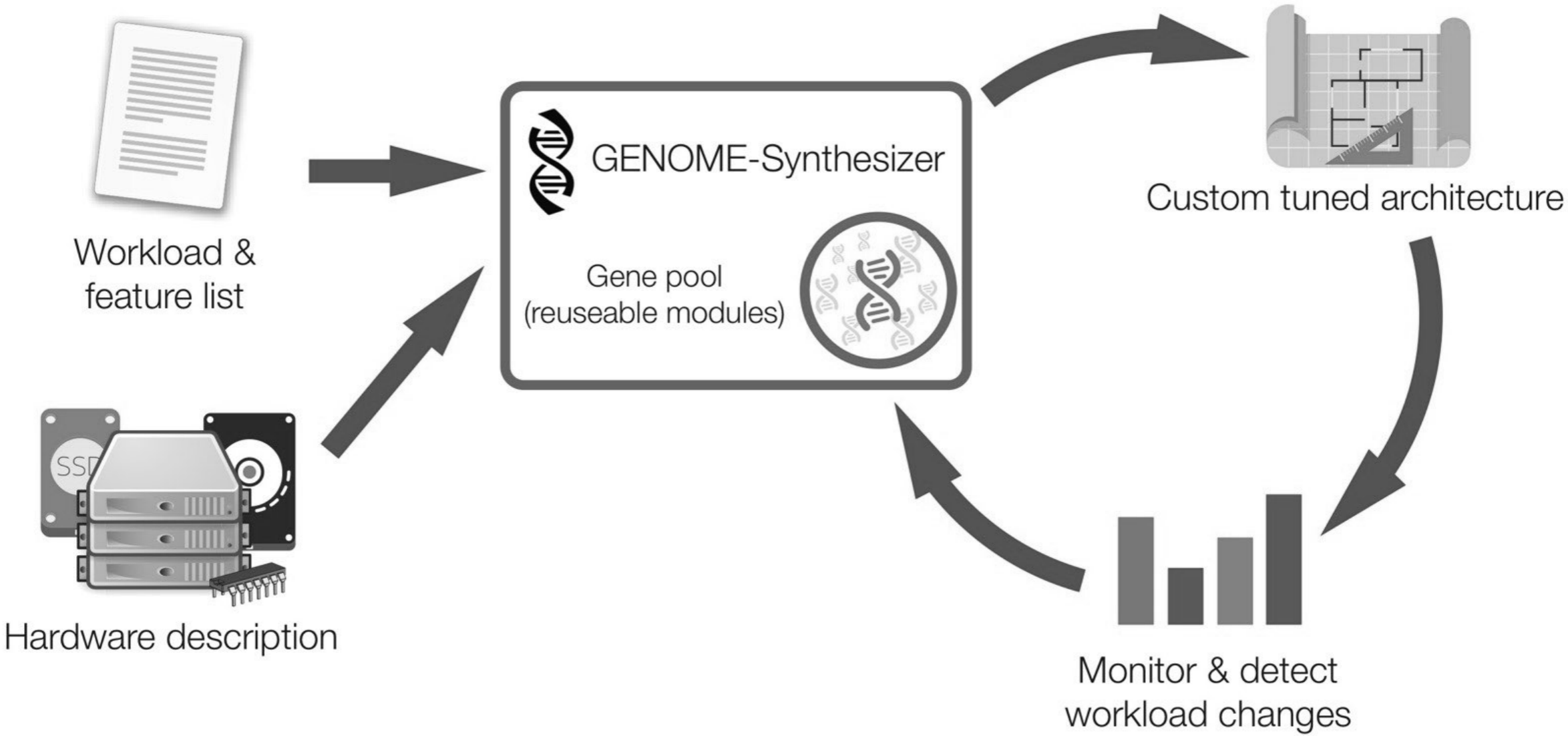
data+queries+hardware



data system

*self-designing
data systems*

adaptivity across architecture borders

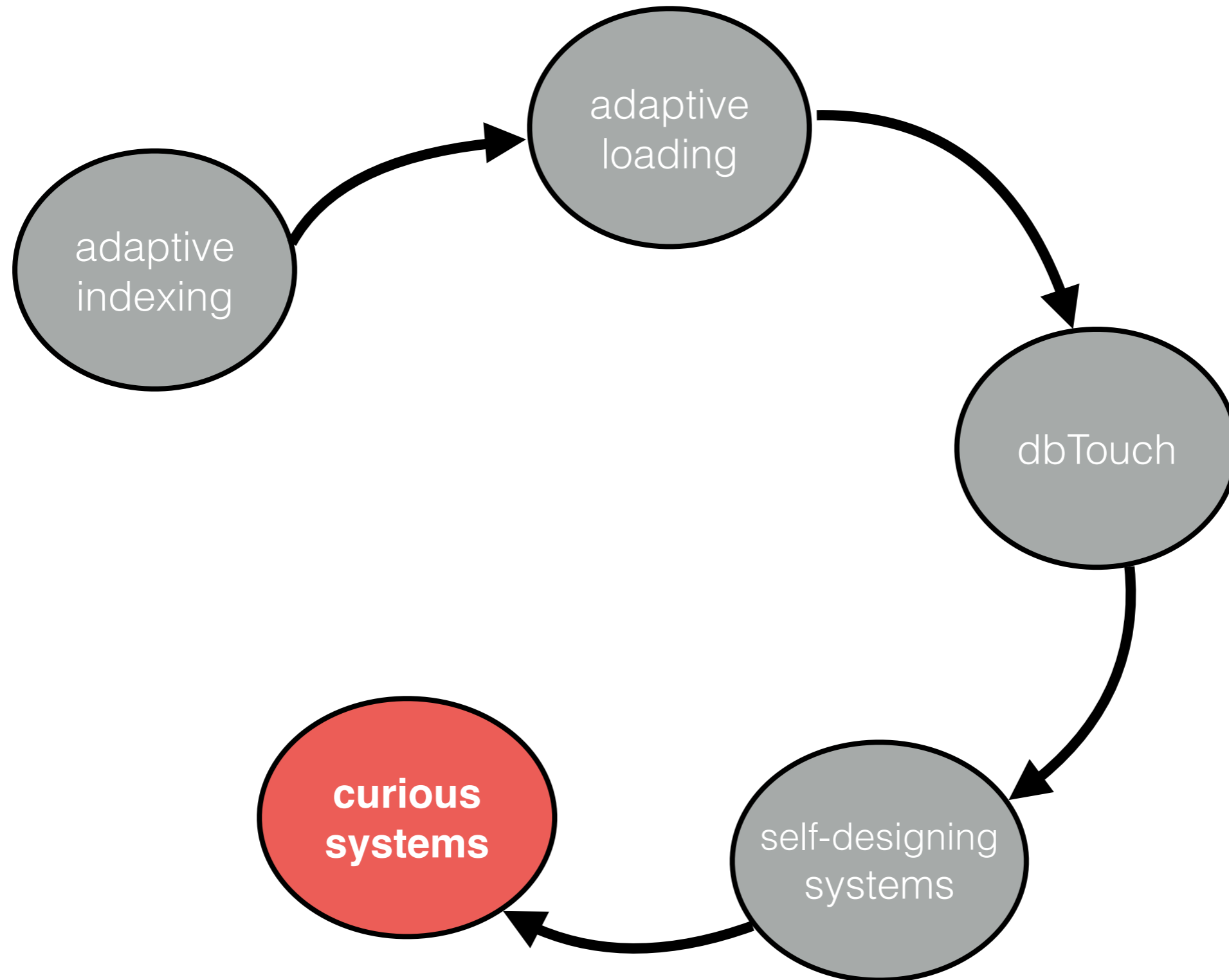


so why self-designing systems?

easy/cheap/fast to design

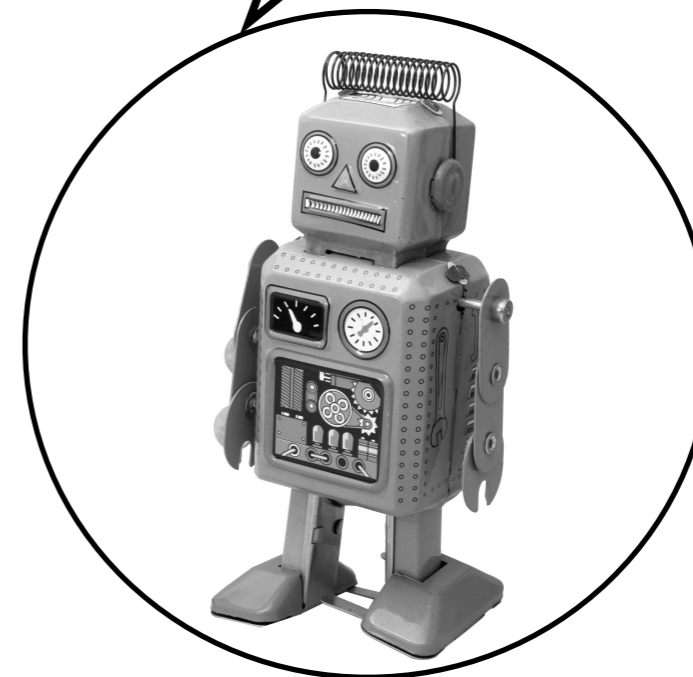
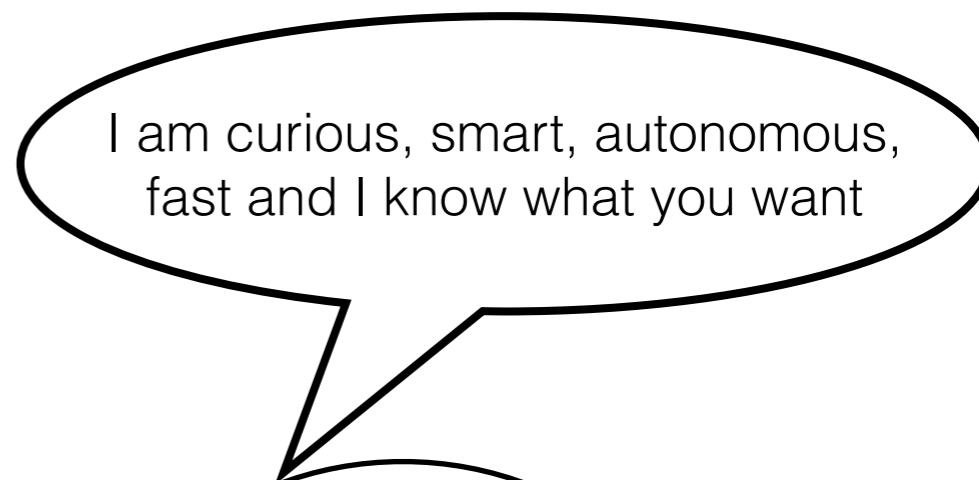
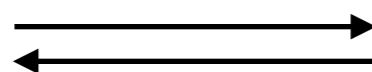
adapt to varying environments

detect suboptimal designs





show me something interesting



DATA

Queriosity

auto-tuning

plug & play

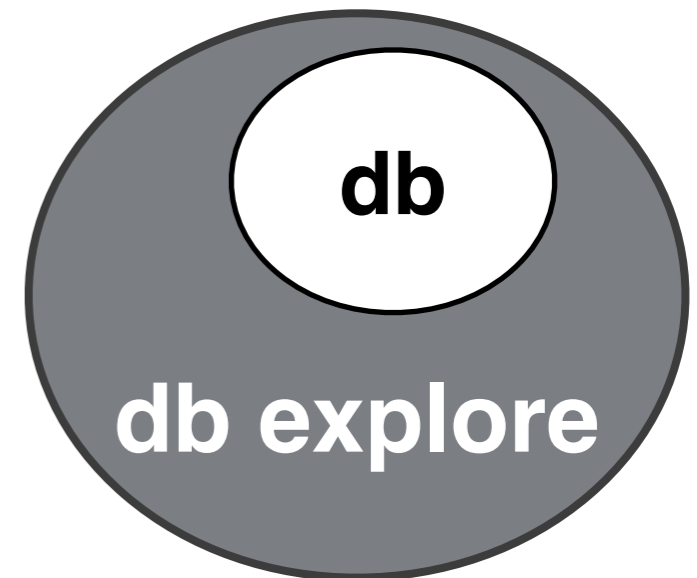
adaptive



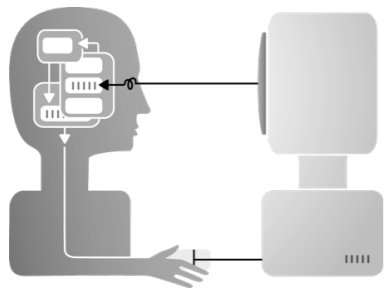
enable data
navigation

interactive

data systems today
allow us to answer queries fast



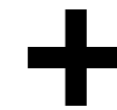
data systems tomorrow
should allow us to find fast which queries to ask



CS, e.g., HCI, OS



data systems



statistics,
sciences, etc.

Joint work with:

Martin Kersten
Stefan Manegold
Goetz Graefe
Harumi Kuno
Anastasia Ailamaki
Ioannis Alagiannis
Miguel Branco
Renata Borovica
Erietta Liarou
Felix Halim
Ronald Yap
Panos Karas
Eleni Petraki
Manos Athanassoulis
Lukas Maas
Abdul Wasay



DATA SYSTEMS. LABORATORY

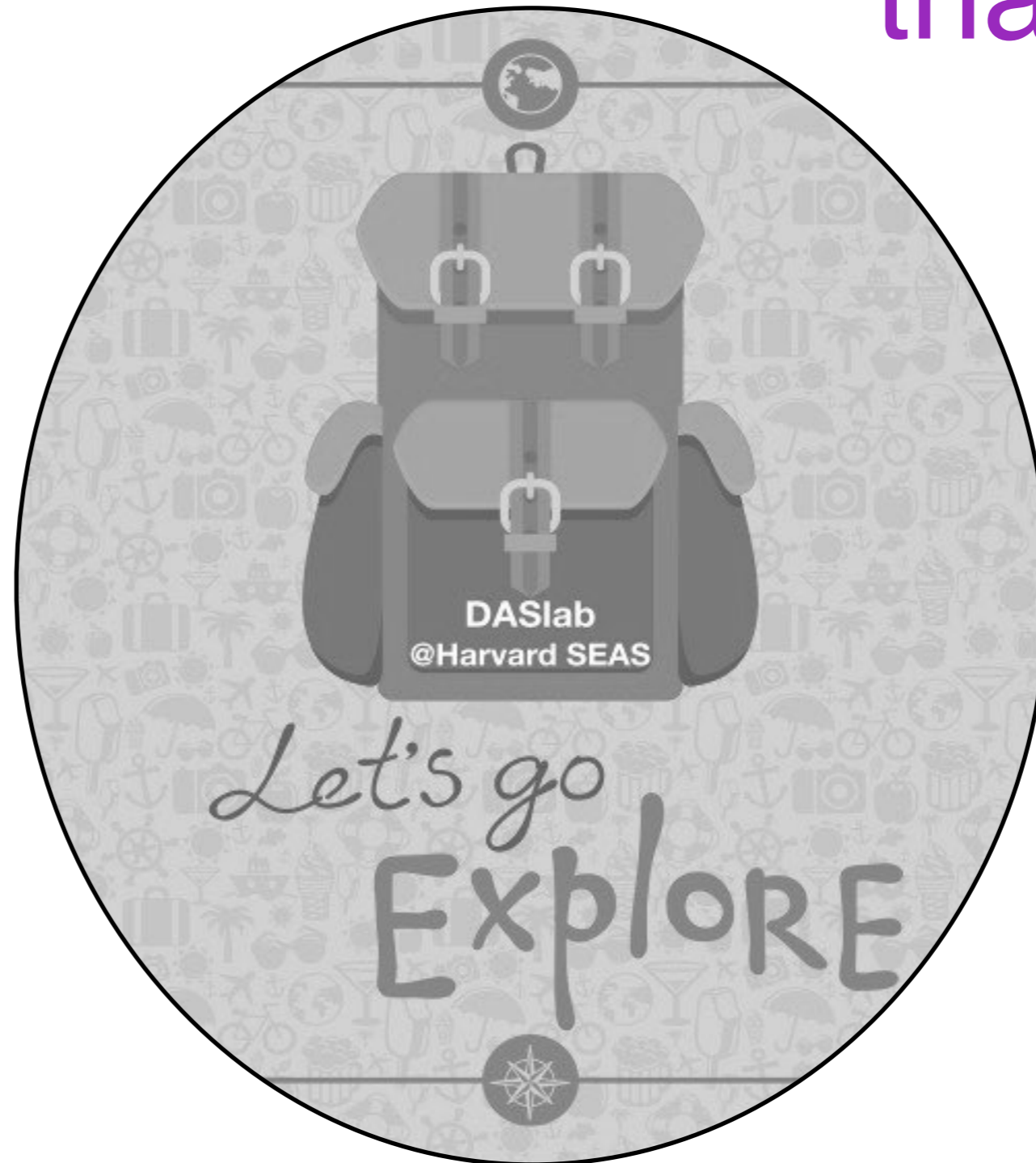
@ Harvard School of Engineering and Applied Sciences

daslab.seas.harvard.edu

thank you!

Joint work with:

Martin Kersten
Stefan Manegold
Goetz Graefe
Harumi Kuno
Anastasia Ailamaki
Ioannis Alagiannis
Miguel Branco
Renata Borovica
Erietta Liarou
Felix Halim
Ronald Yap
Panos Karas
Eleni Petraki
Manos Athanassoulis
Lukas Maas
Abdul Wasay



DATA SYSTEMS. LABORATORY
@ Harvard School of Engineering and Applied Sciences

daslab.seas.harvard.edu