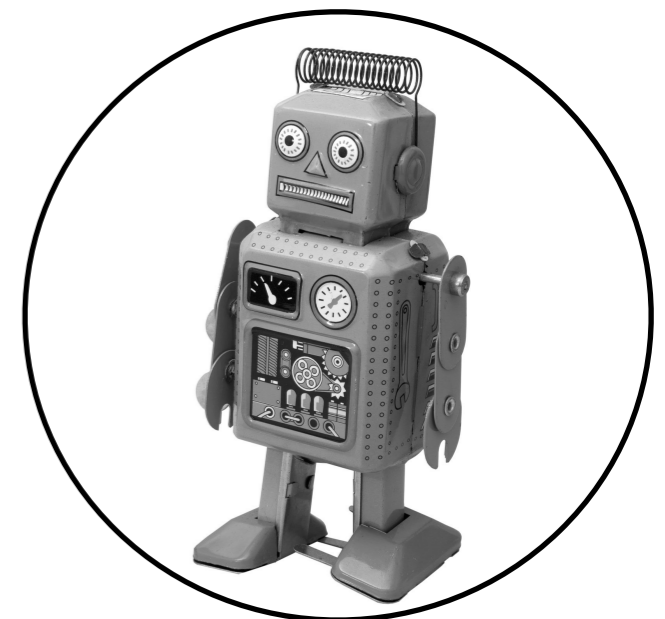
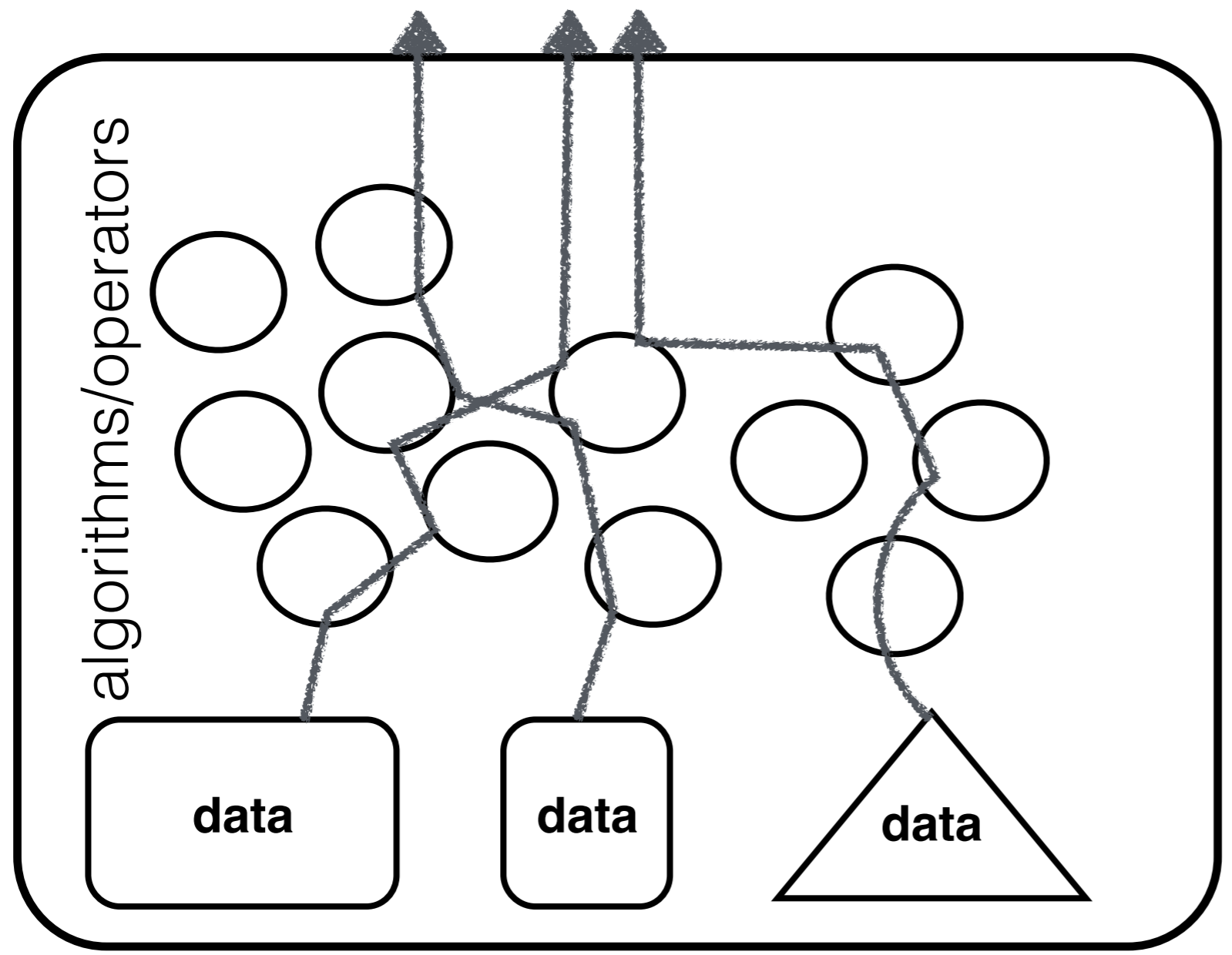
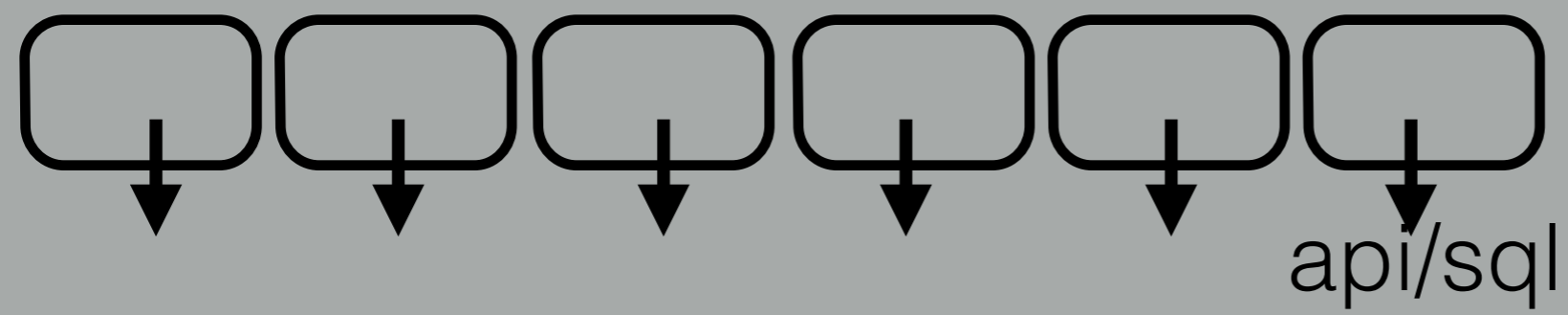


Data Systems that are Easy to Design, Tune and Use

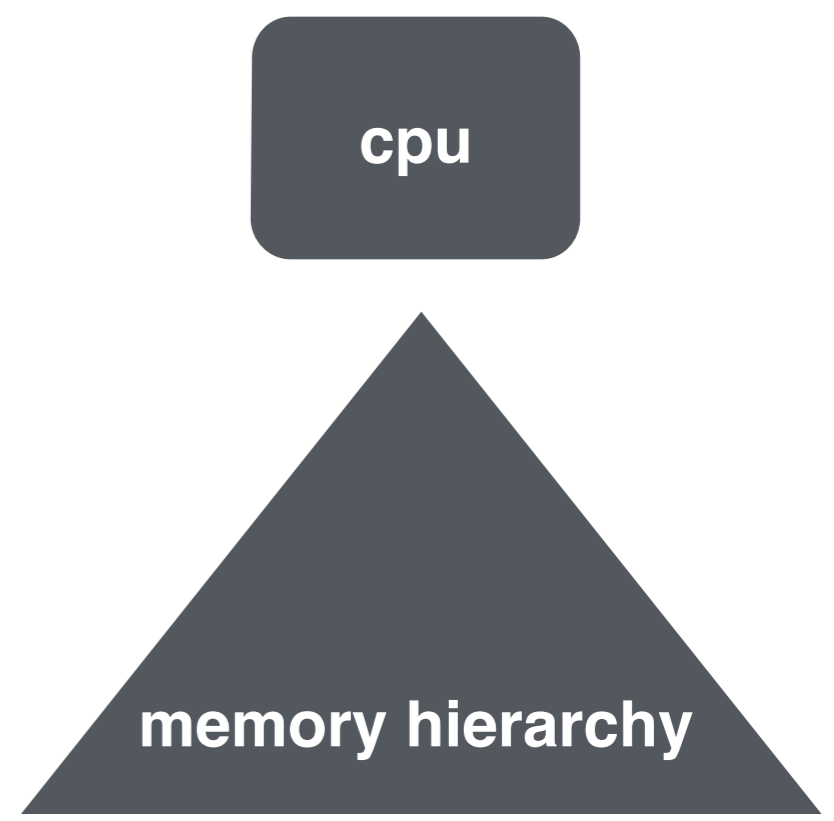
Stratos Idreos

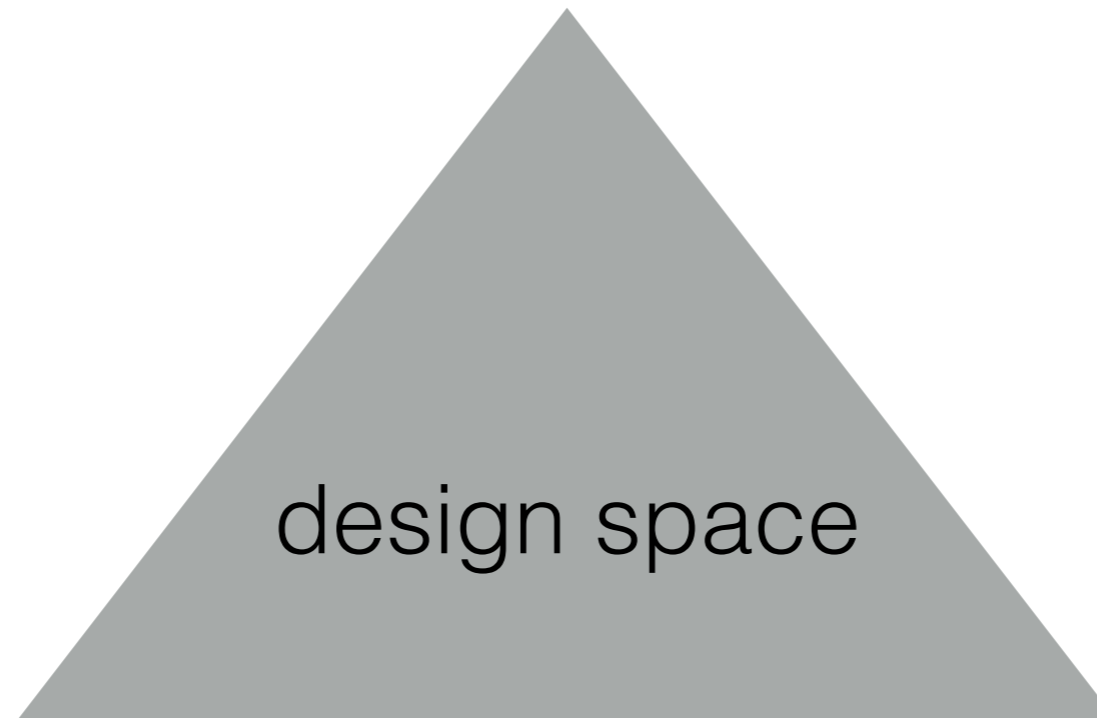


applications

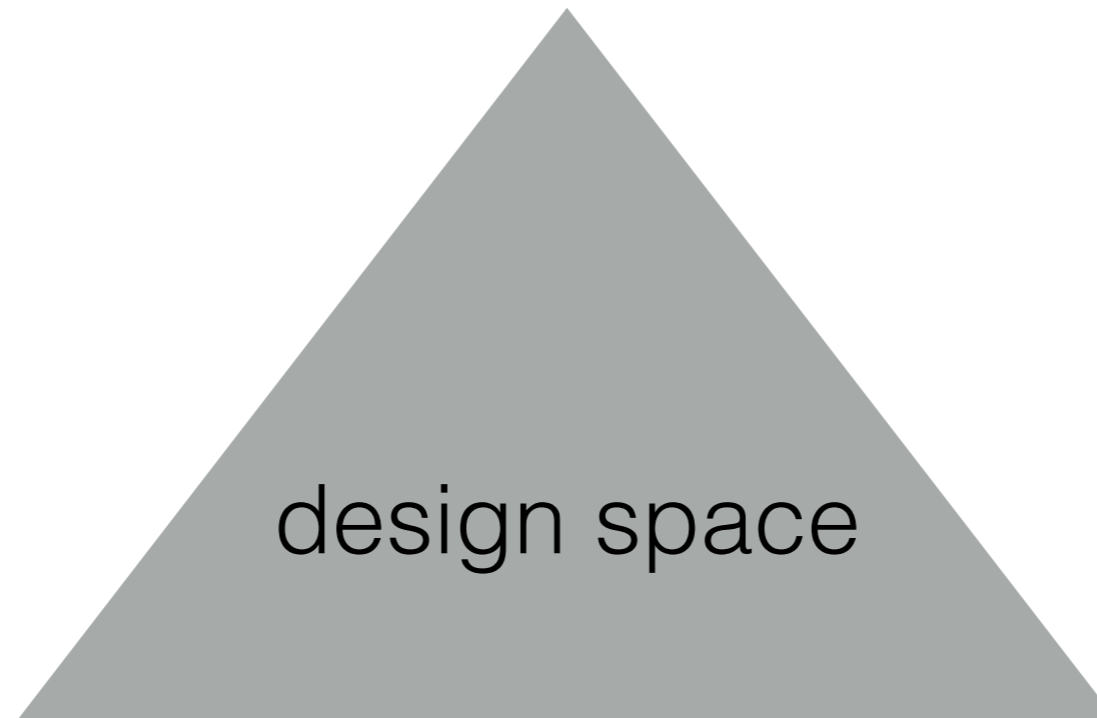


data system kernel





it all starts with how we store data
every bit matters



it all starts with how we store data
every bit matters

no fixed decisions
from static to dynamic designs

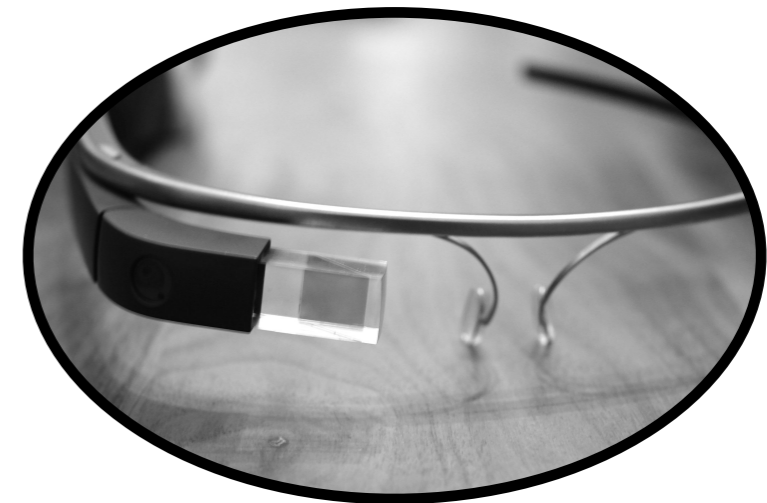
today



today



tomorrow



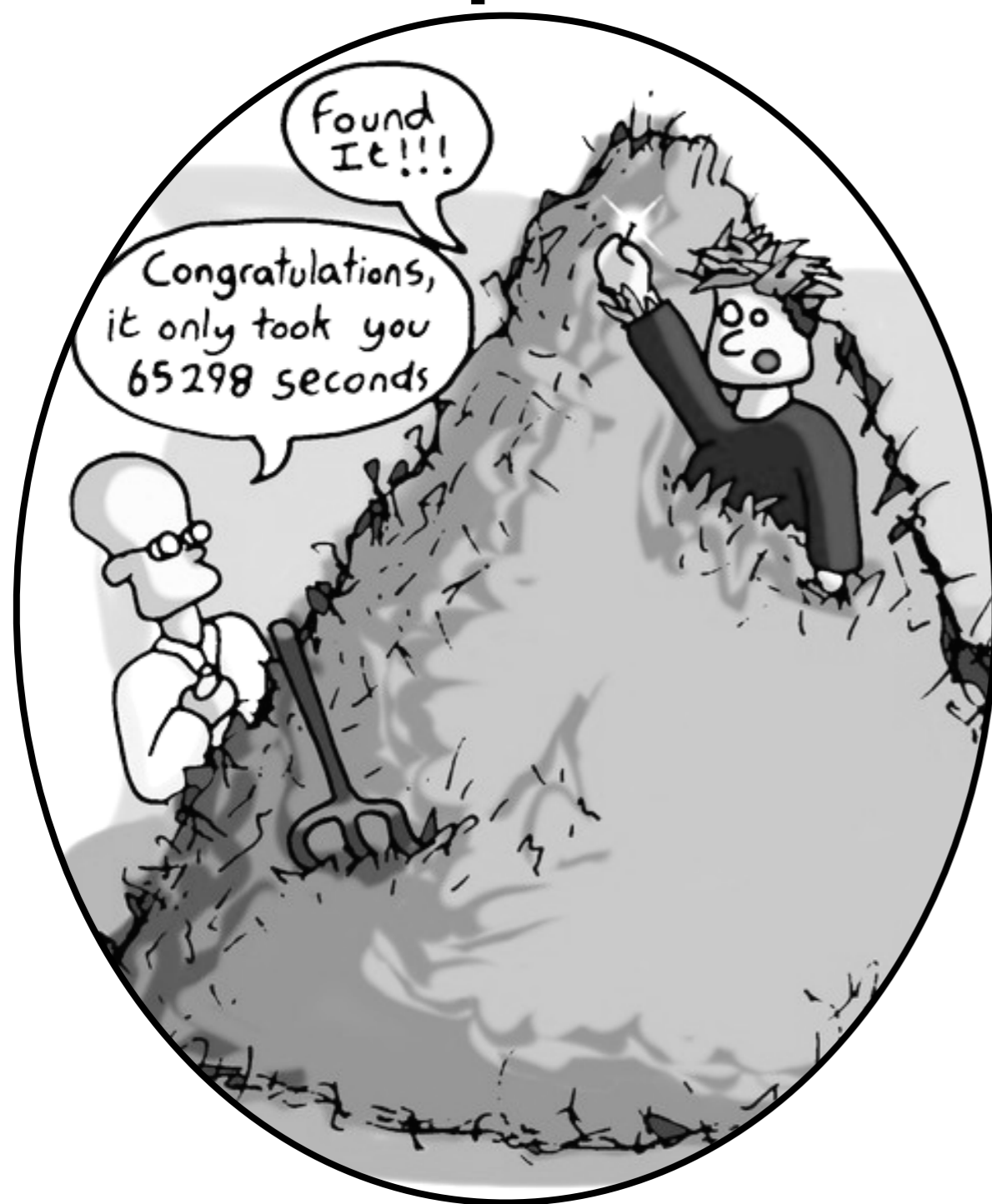
soon everyone will need to be a “data scientist”

**hmm,
my data is too big :(**





data exploration



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

data has always been **big**

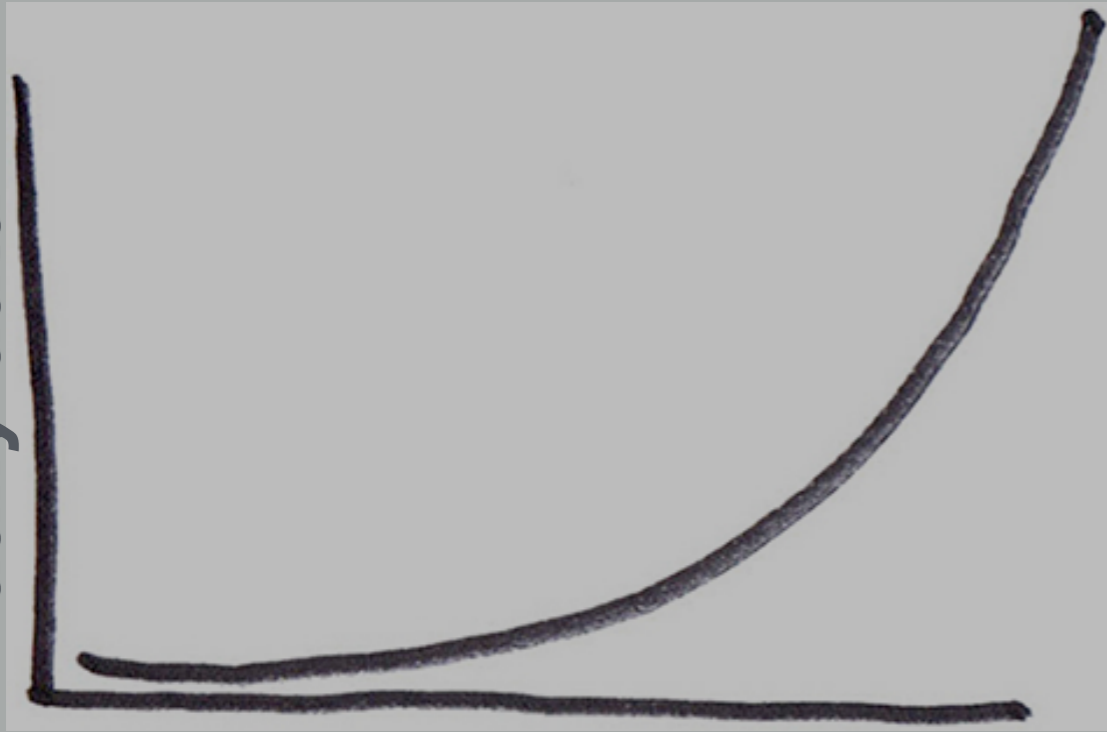
volume

velocity

variety

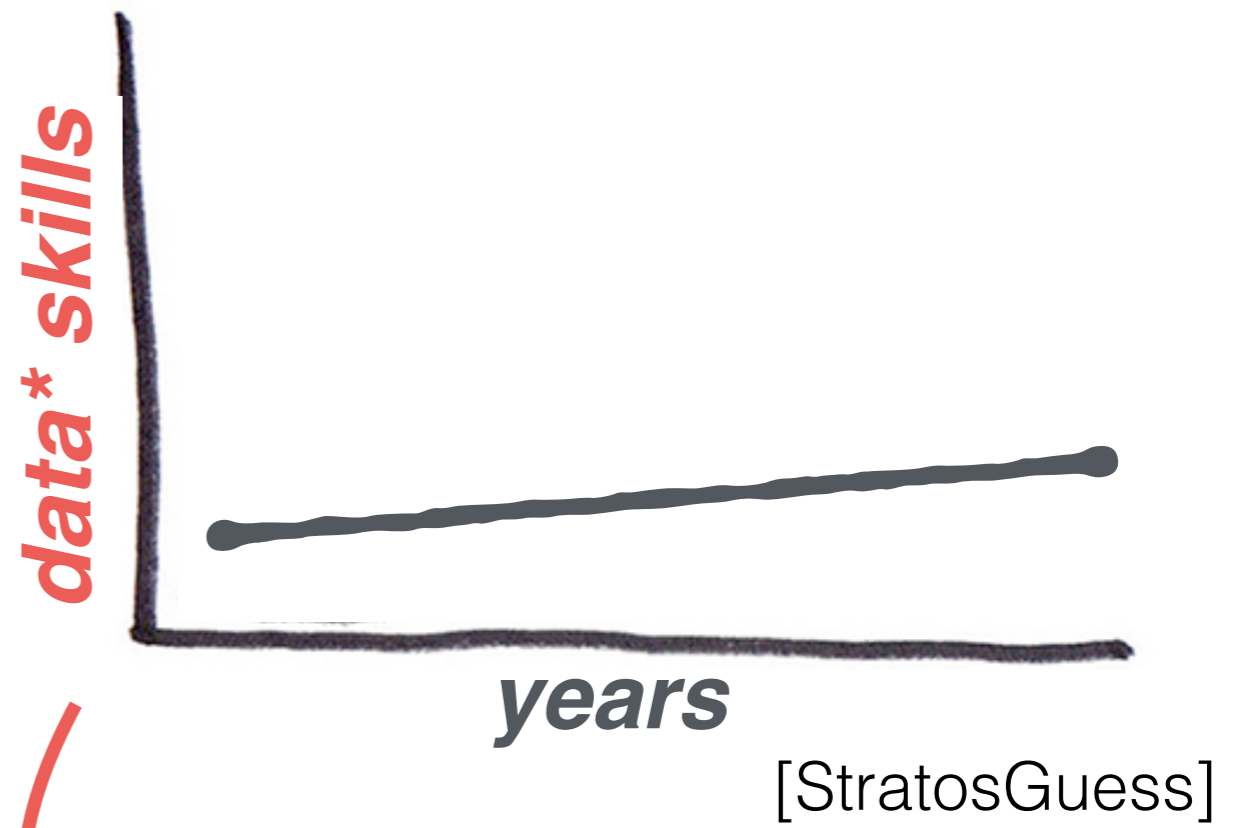
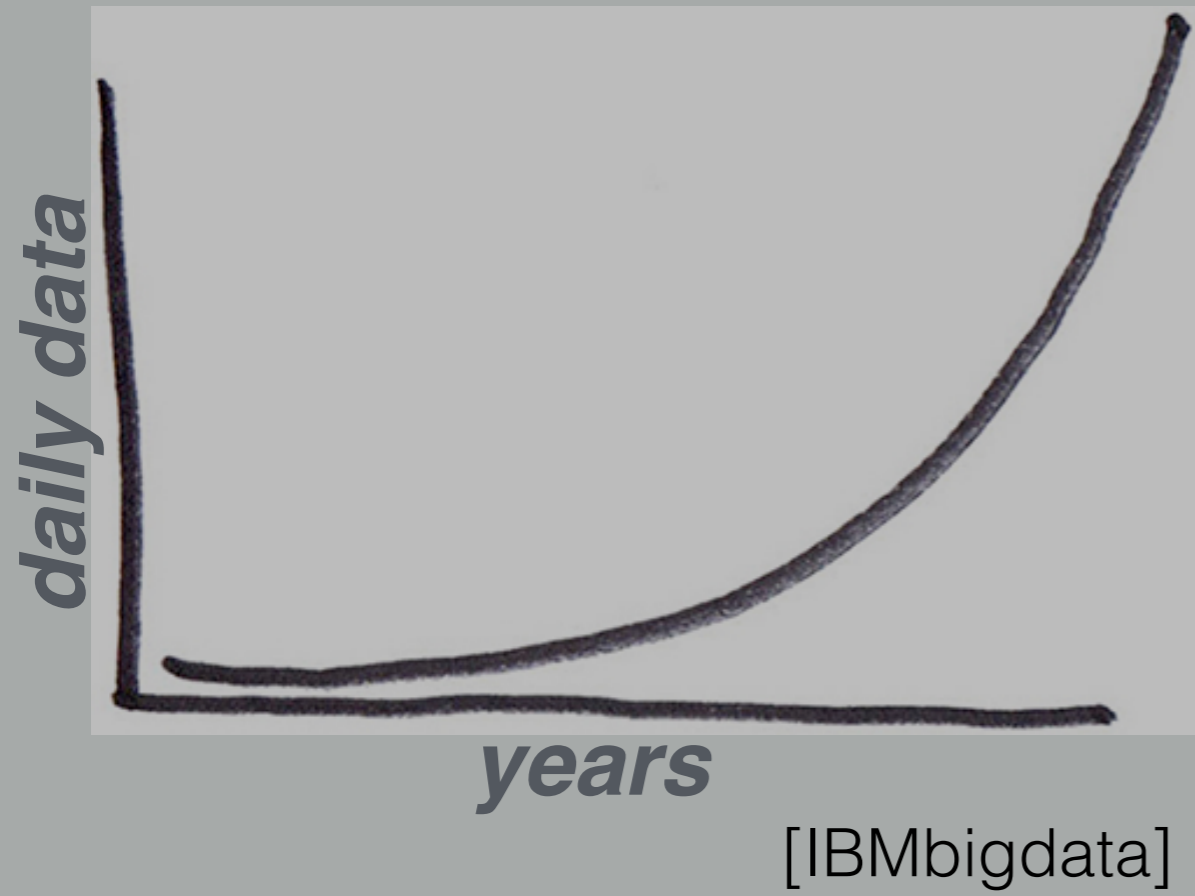
veracity

daily data



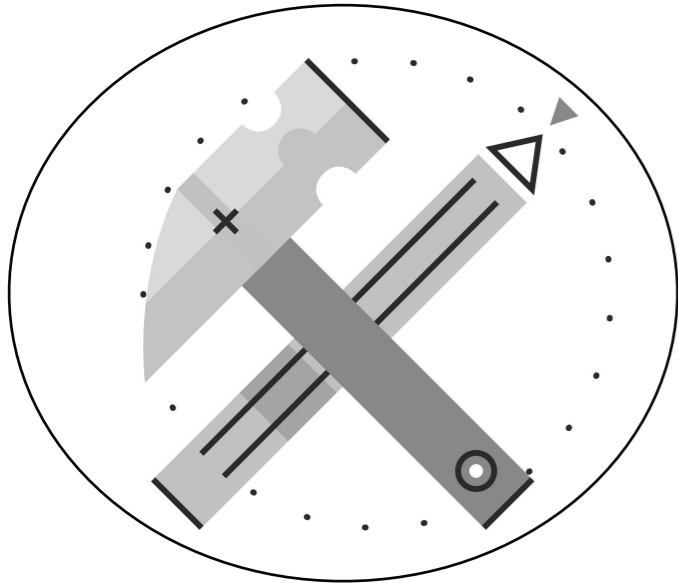
years

[IBMbigdata]



data system design,
set-up, tune, use

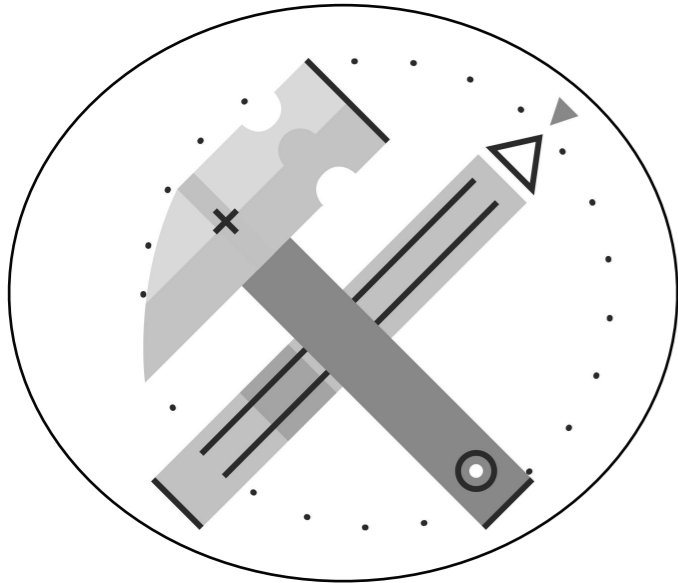
data systems that are **easy** to:



(years)

design & build

data systems that are **easy** to:



(years)

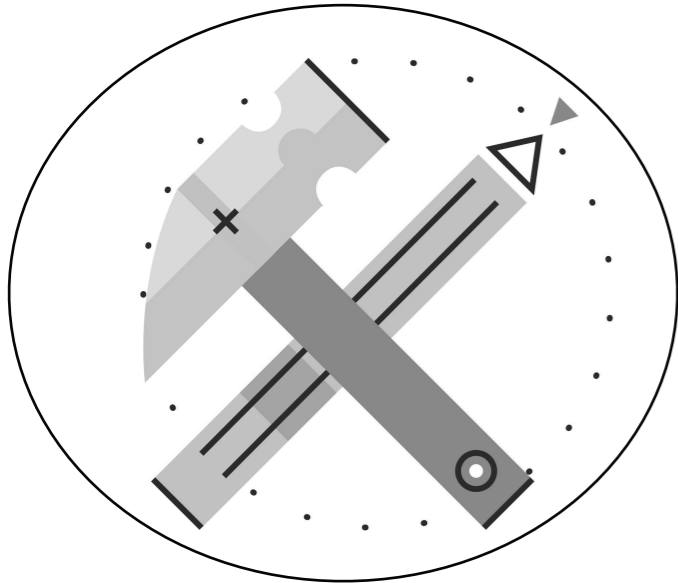
design & build



(months)

set-up & tune

data systems that are **easy** to:



(years)

design & build



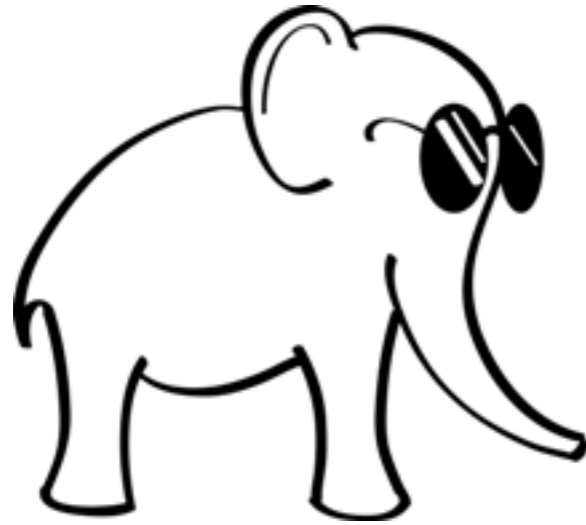
(months)

set-up & tune

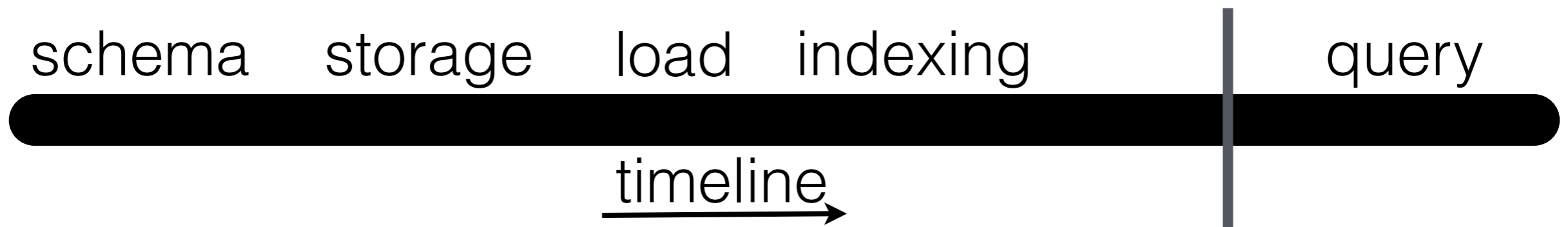


(hours/days)

use



too many preparation options
lead to complex installation



expert users - idle time - workload knowledge



users/applications
declarative interface
ask what you want



 **DBA**

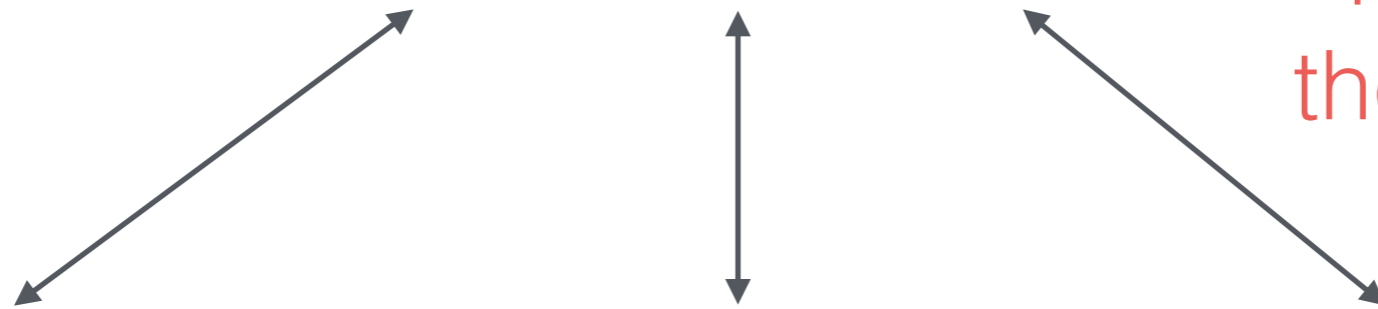


db system



users/applications

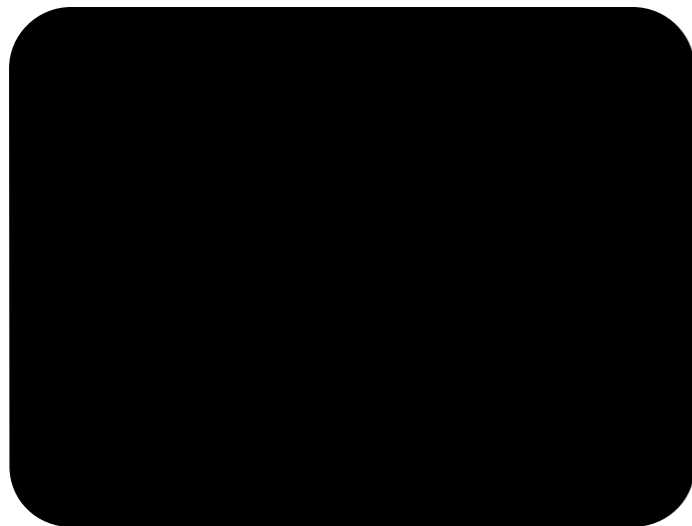
need to choose
the proper system
&
workloads/
applications
change rapidly



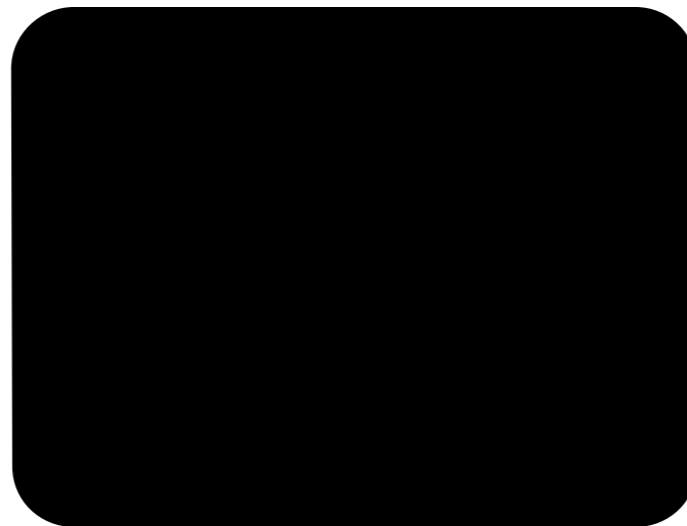
DBA1



DBA2



data system 1

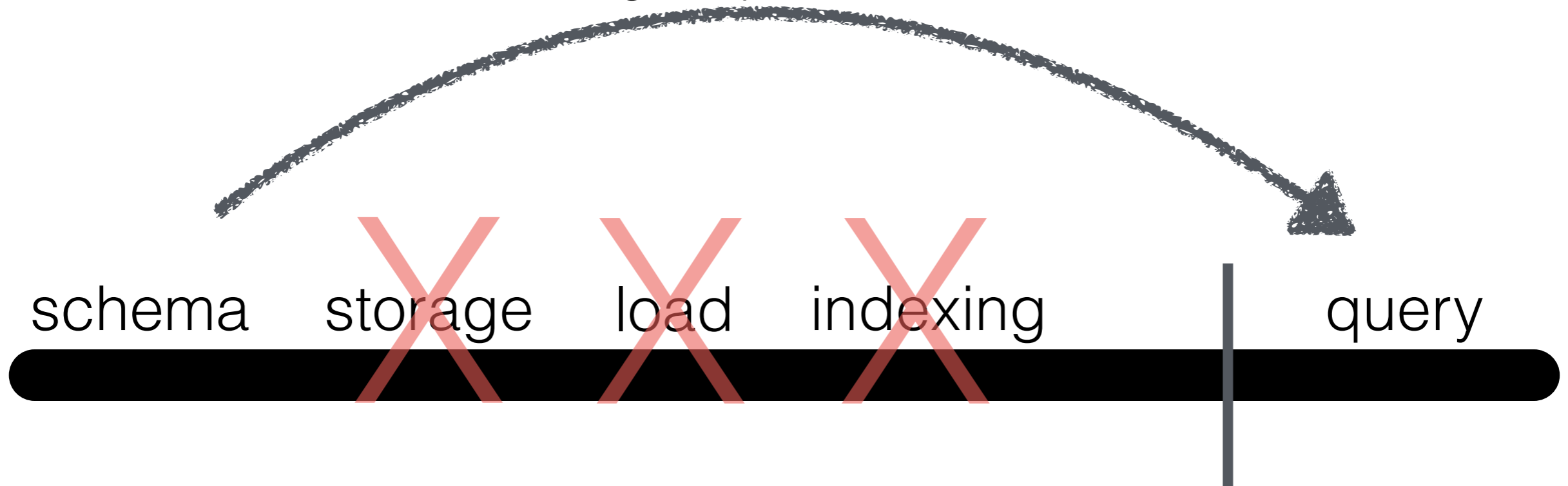


data system 2

...

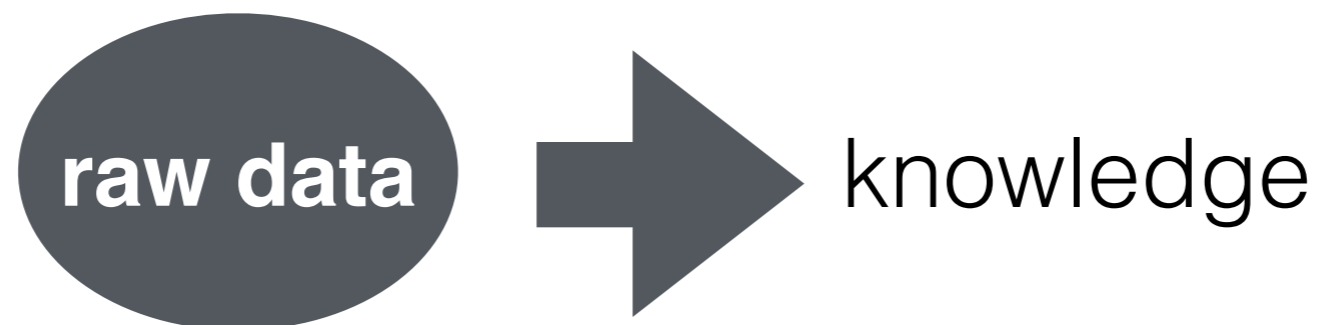
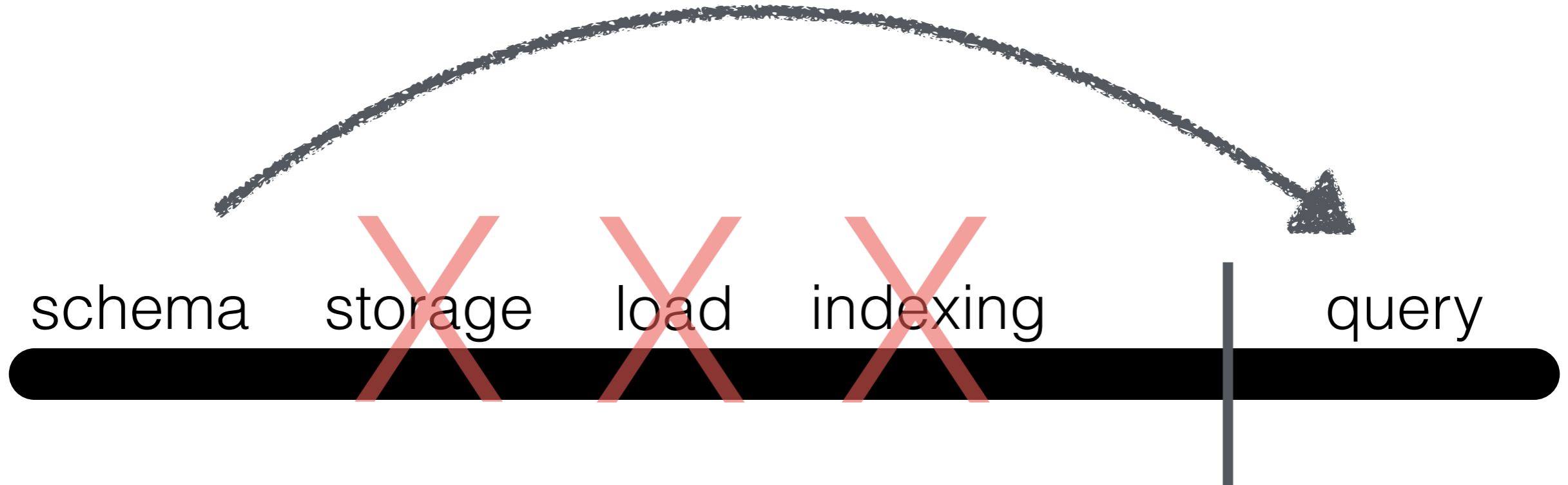


be able to query the data immediately
& with good performance





be able to query the data immediately
& with good performance



indexing



tune= create proper indices offline
performance 10-100X

indexing

storage

load

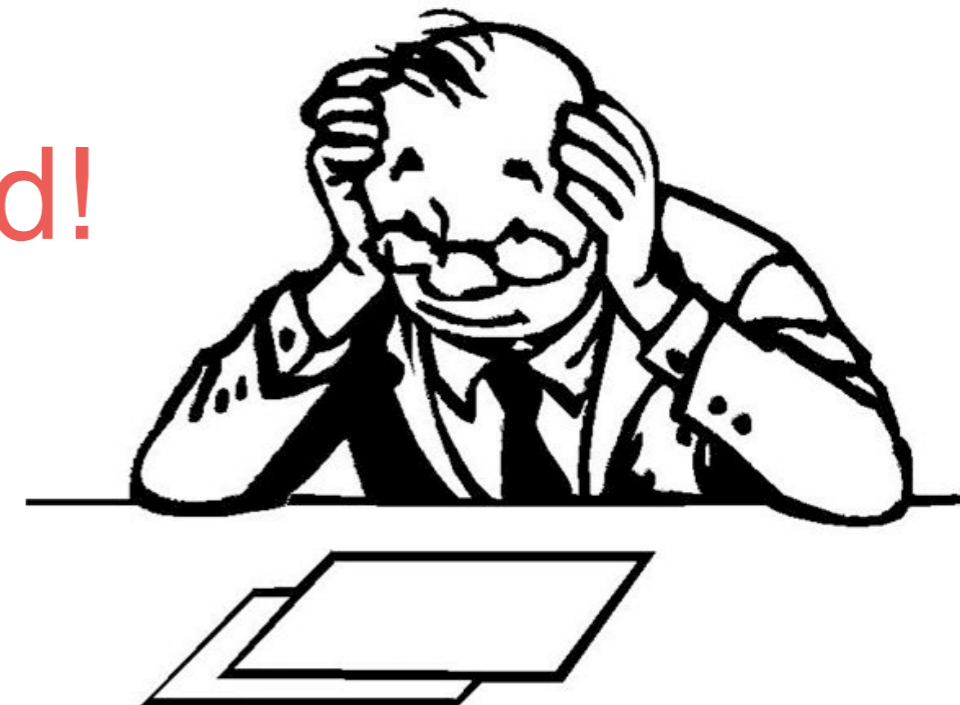
indexing

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?



storage

load

indexing

query

storage

load

indexing

query

timeline

storage

load

indexing

query

sample workload

timeline

storage

load

indexing

query

sample workload

analyze

timeline

storage

load

indexing

query

sample workload

analyze

create indices

timeline



storage

load

indexing

query

sample workload

analyze

create indices

query

timeline



storage

load

indexing

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

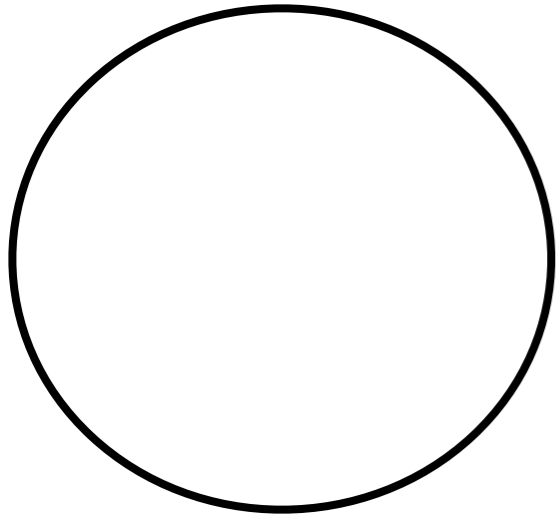
~~idle time~~

~~workload
knowledge~~

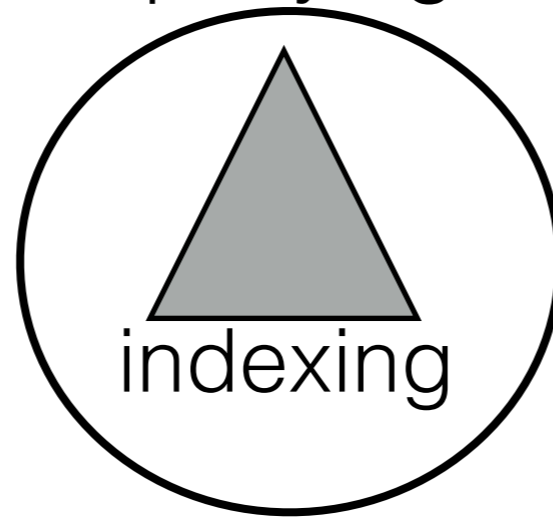
~~external
tools~~

~~human
control~~

initialization



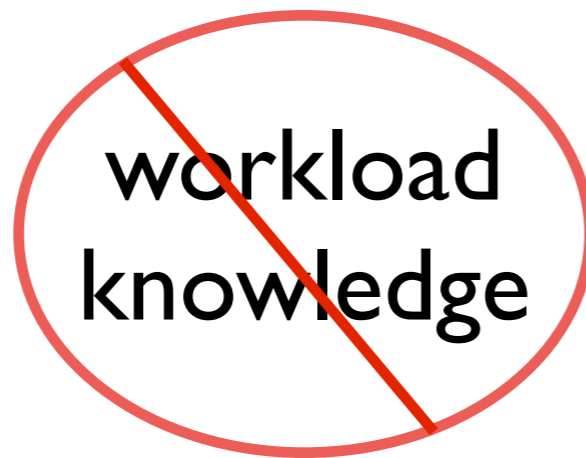
querying



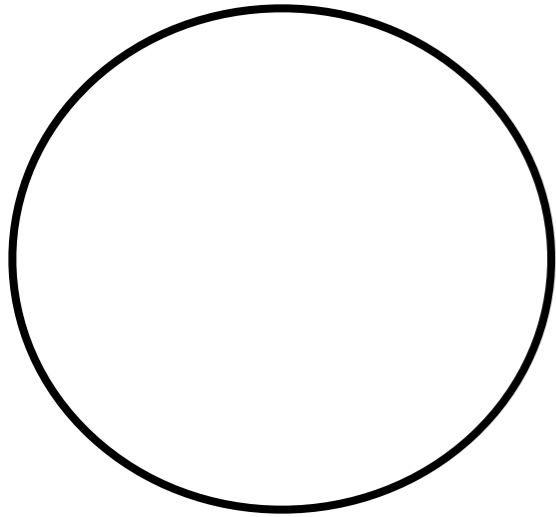
database cracking

auto-tuning database kernels

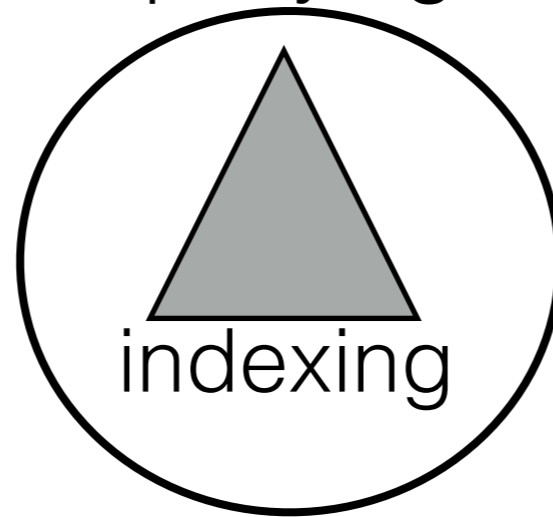
incremental, adaptive, partial indexing



initialization



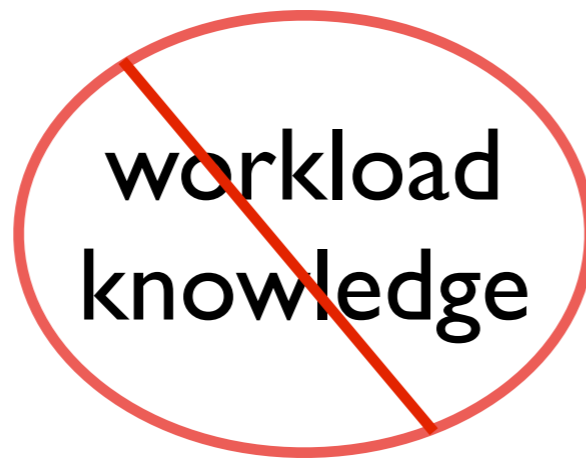
querying



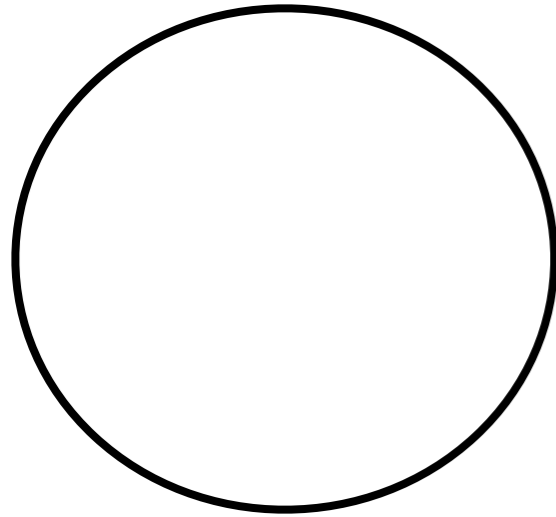
database cracking

auto-tuning database kernels

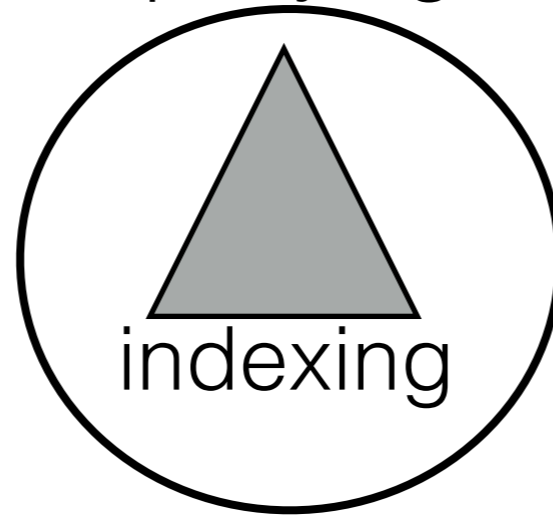
incremental, adaptive, partial indexing



initialization



querying



database cracking

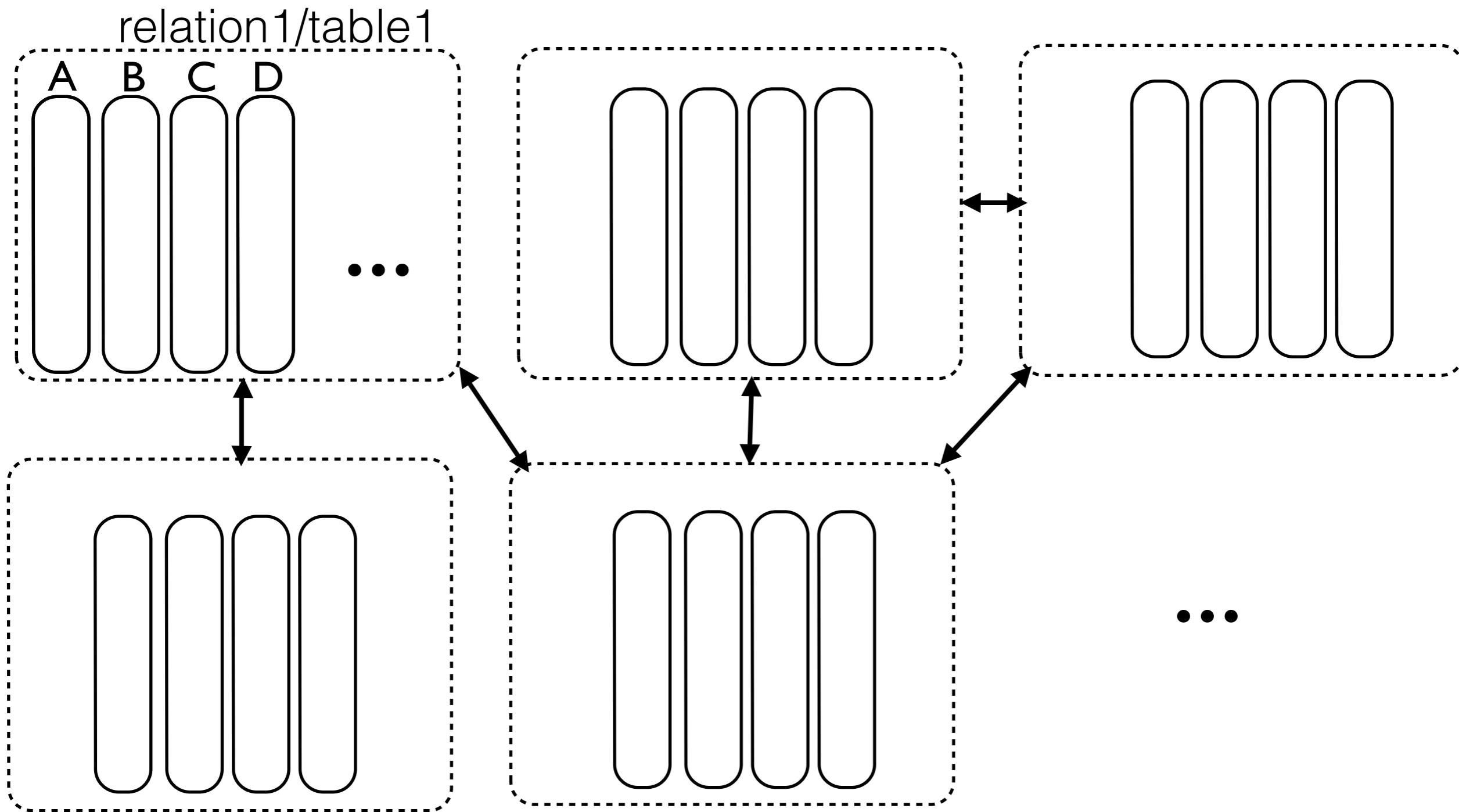
auto-tuning database kernels

incremental, adaptive, partial indexing

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

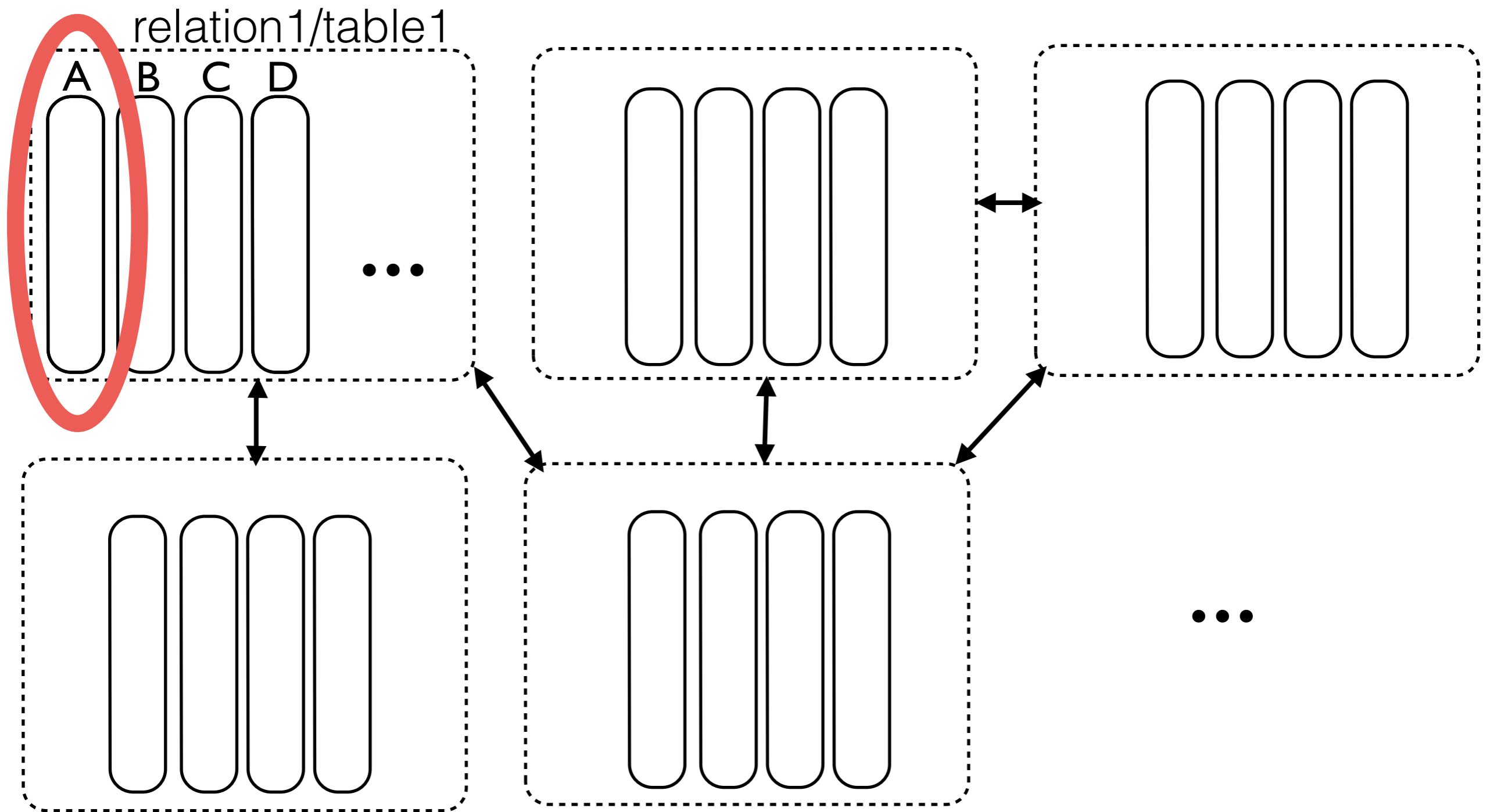
column-store database

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

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



- 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

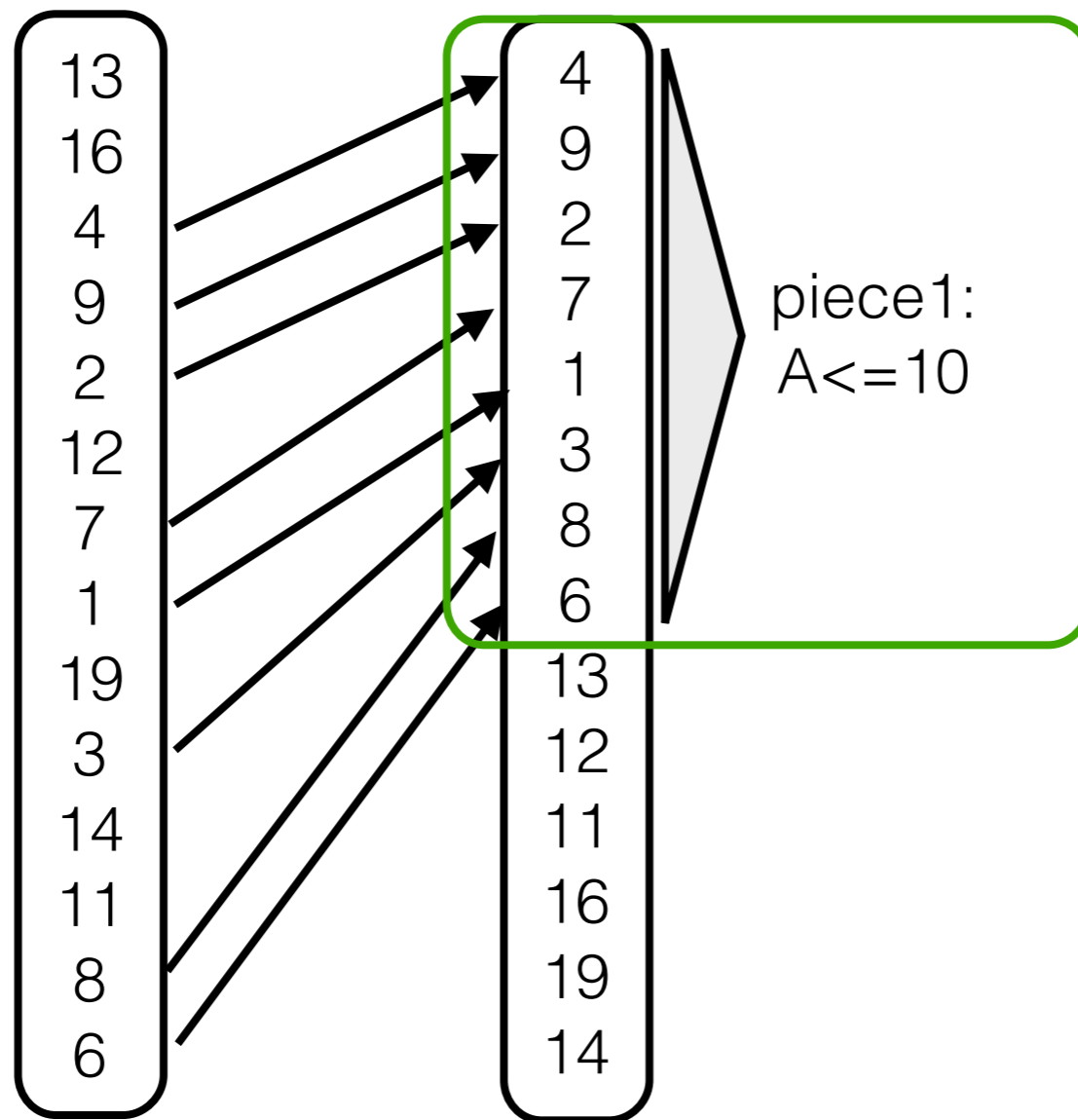
column A

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

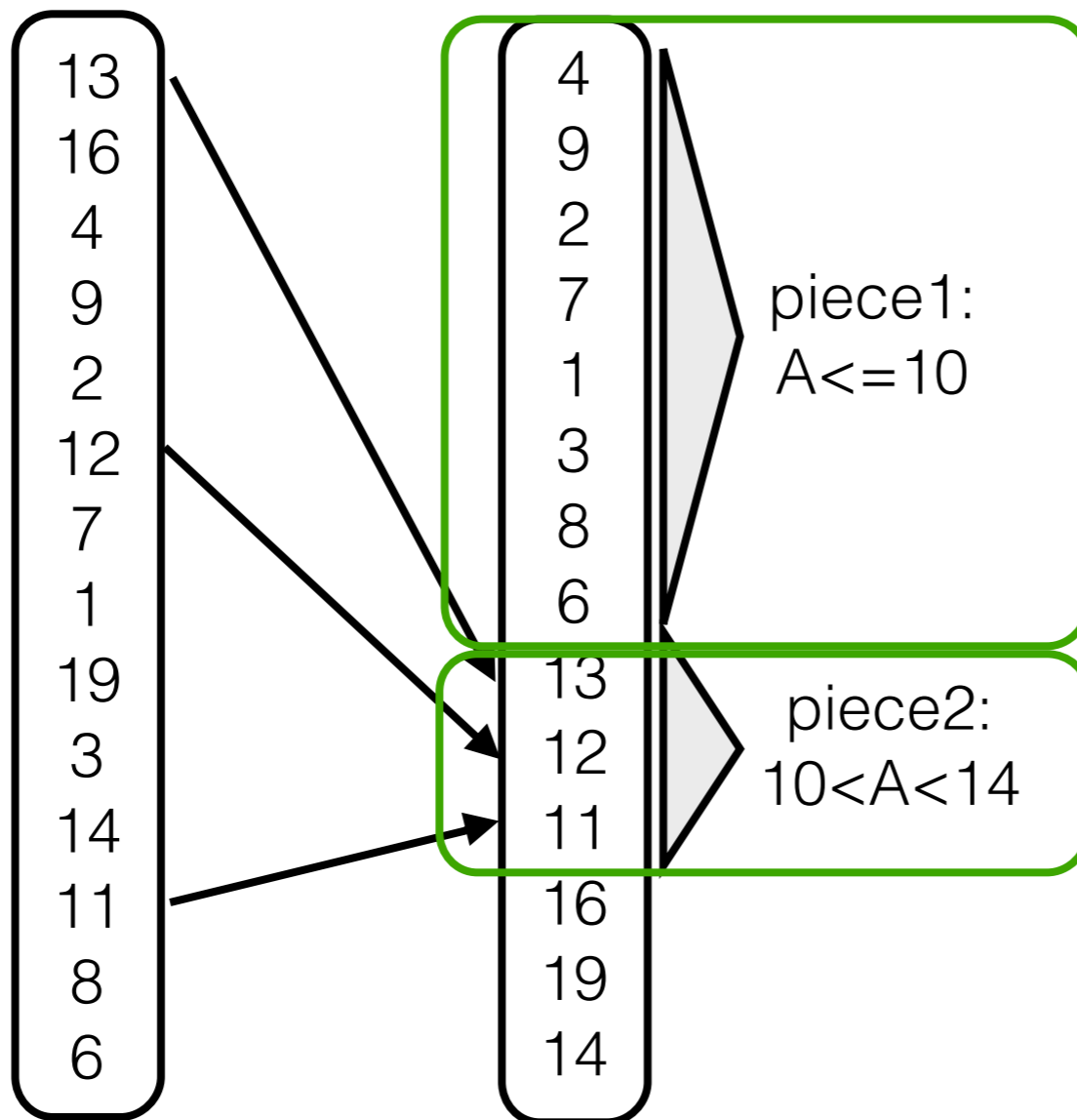
Q1:
select R.A
from R
where R.A > 10
and R.A < 14

column A



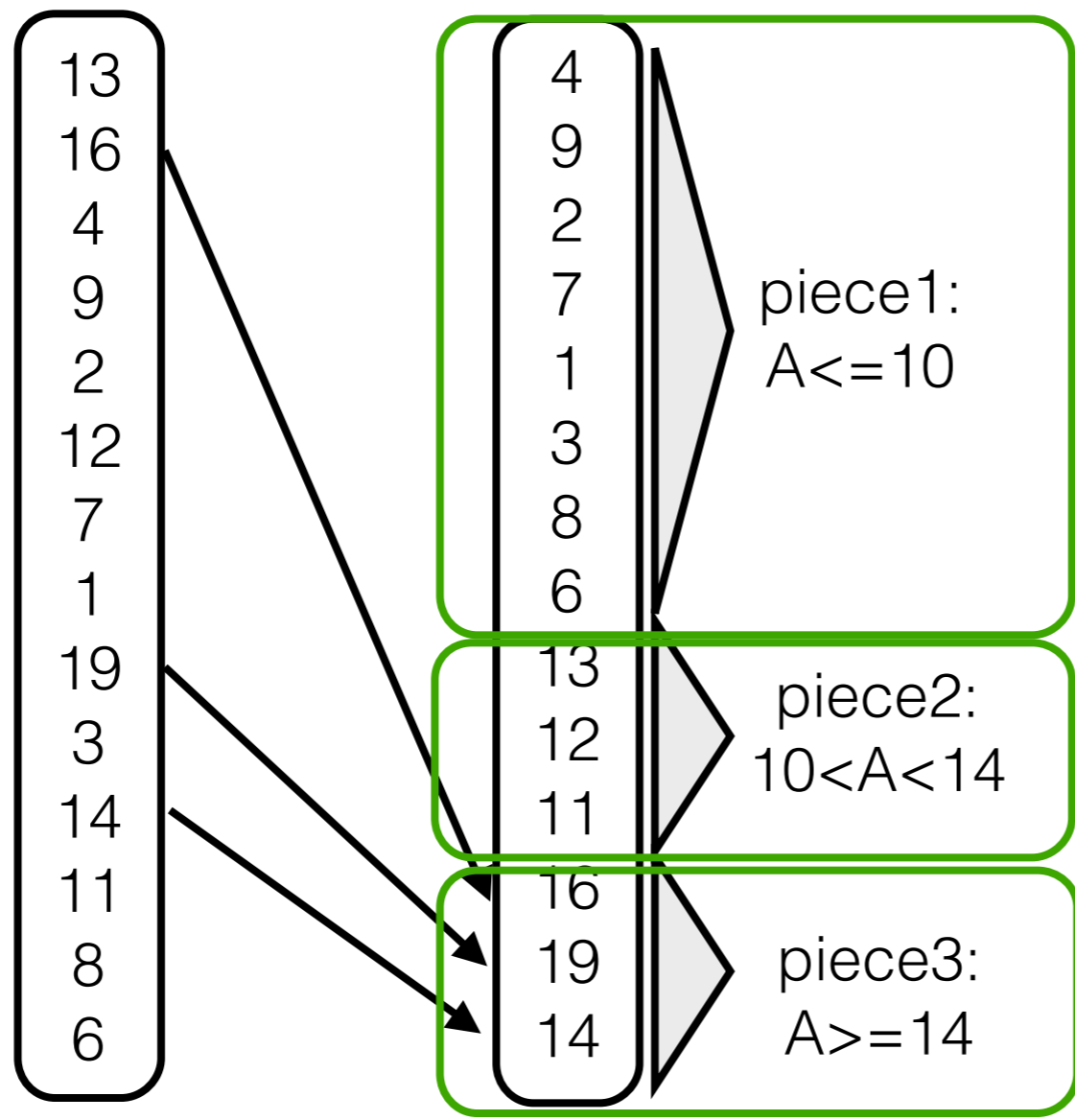
Q1:
select R.A
from R
where R.A > 10
and R.A < 14

column A



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

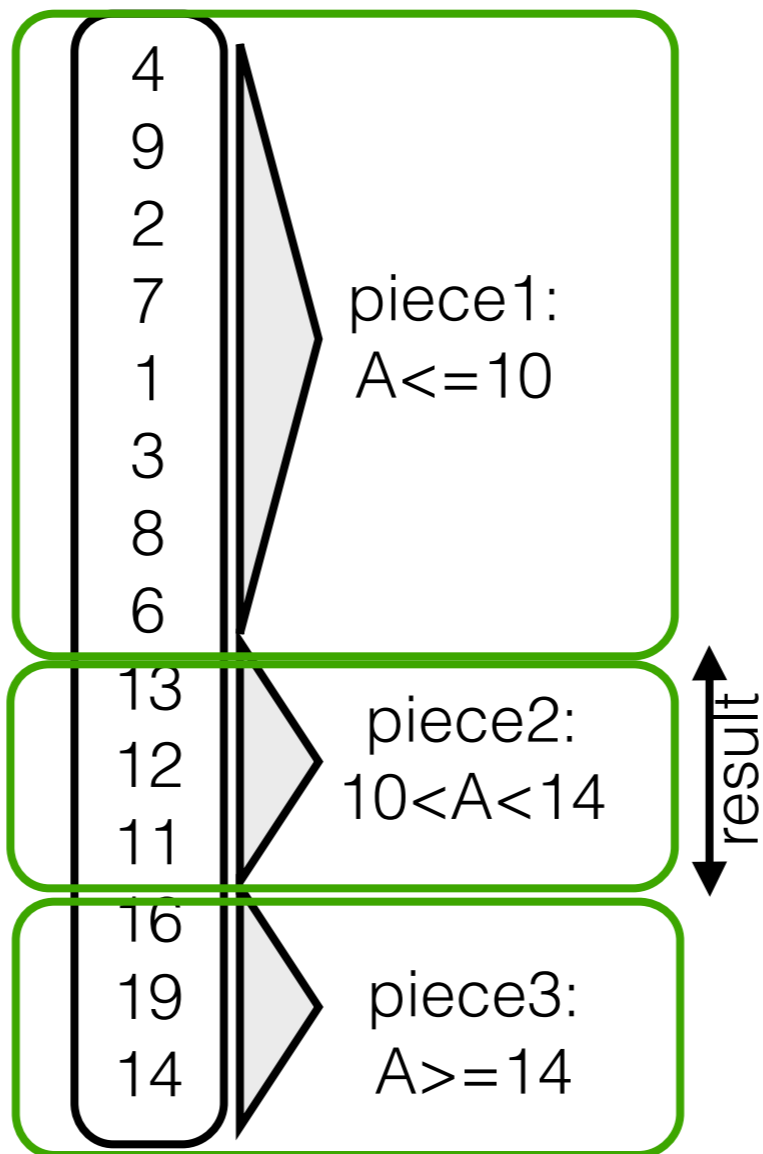
column A



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

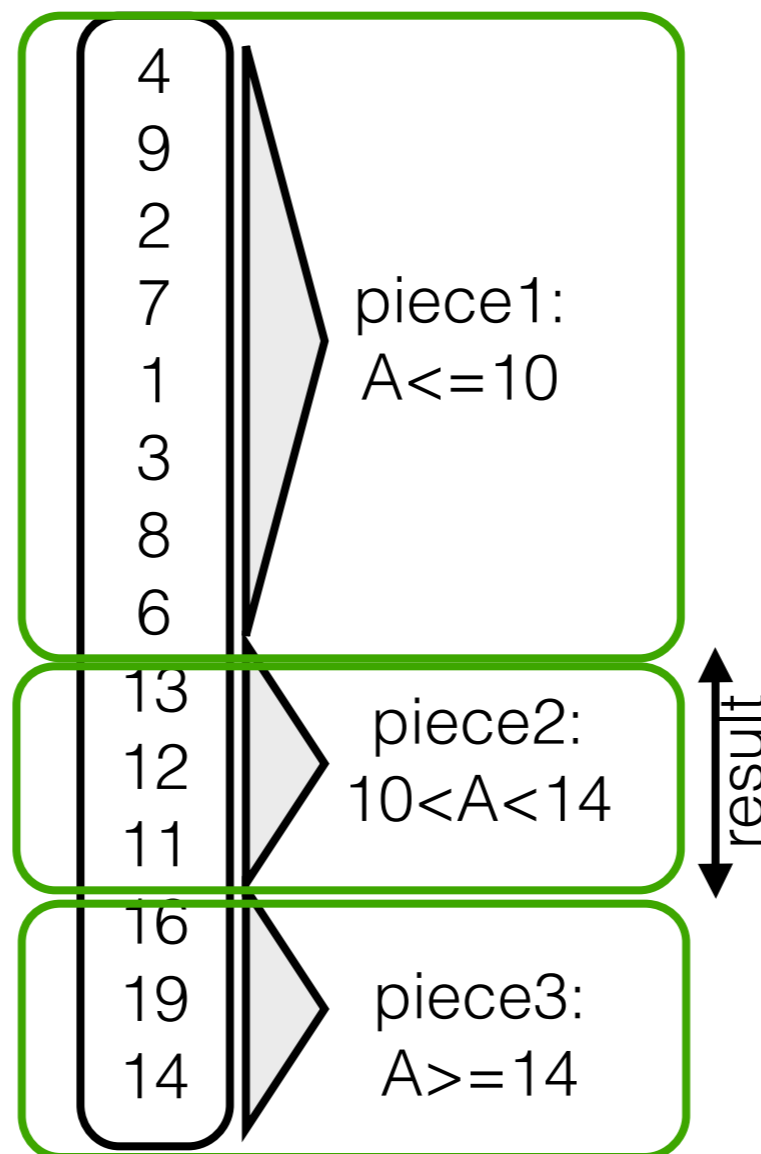


gain knowledge on how data is organized

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

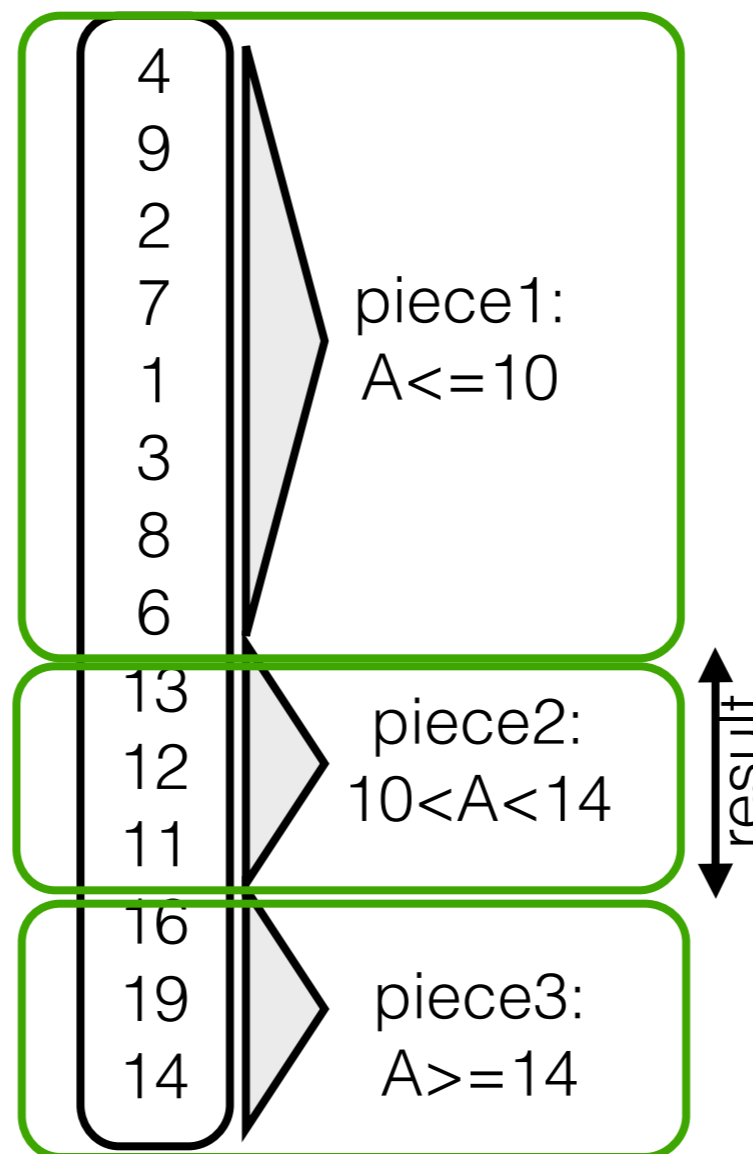


gain knowledge on how data is organized

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



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

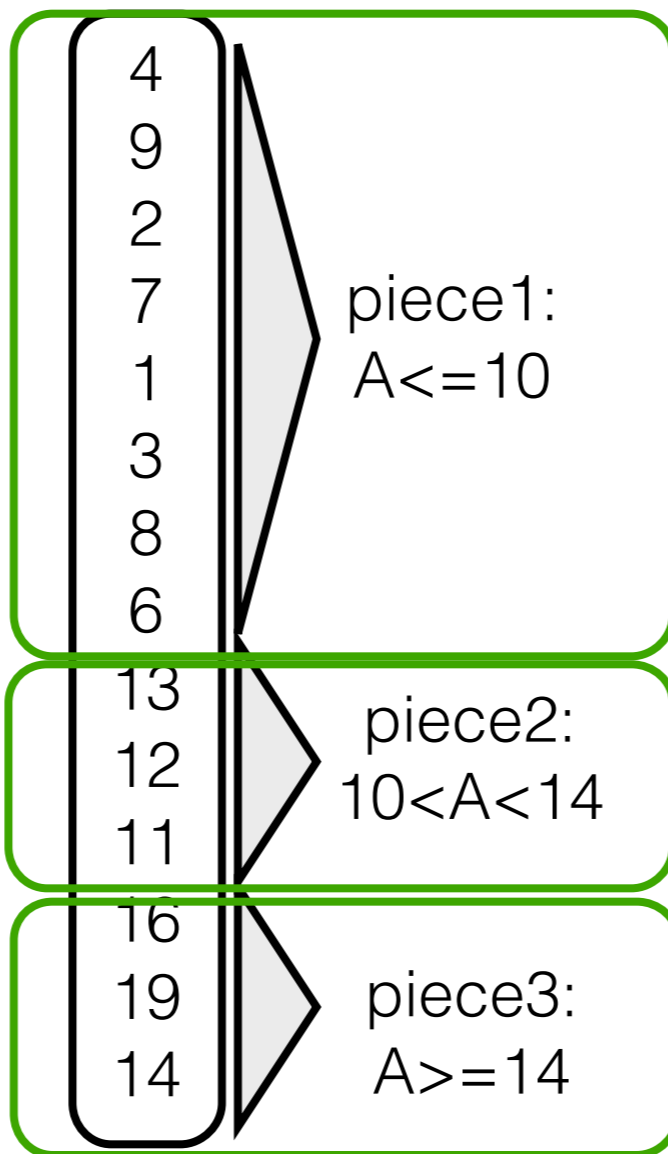
Database Cracking CIDR 2007

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



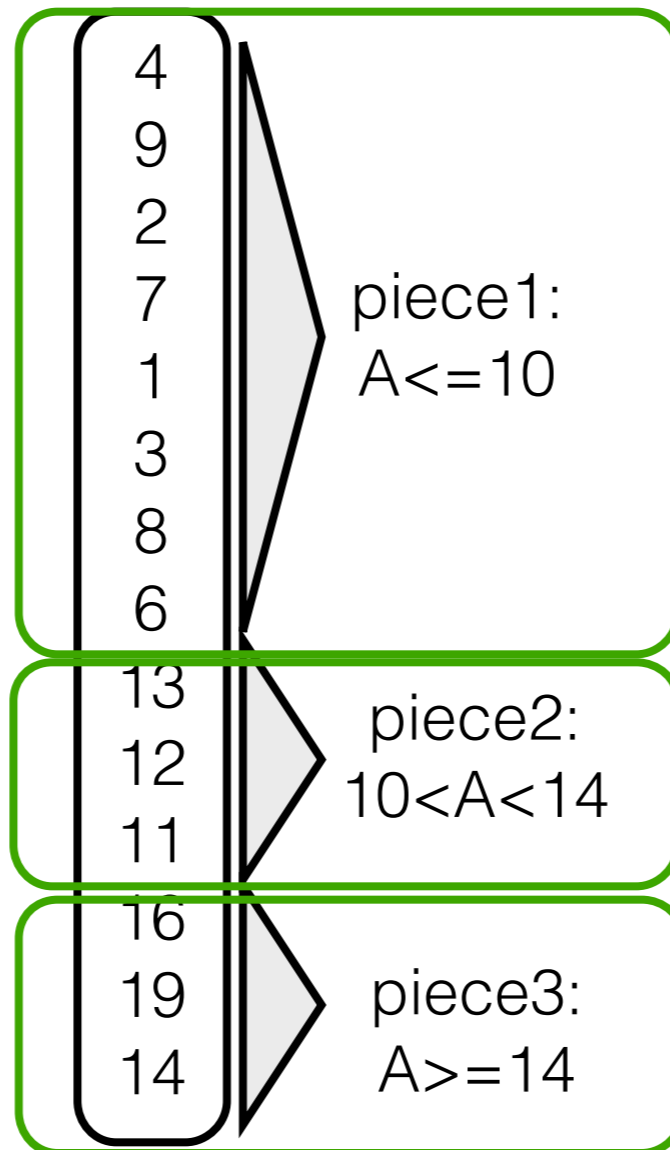
dynamically/on-the-fly within the select-operator

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

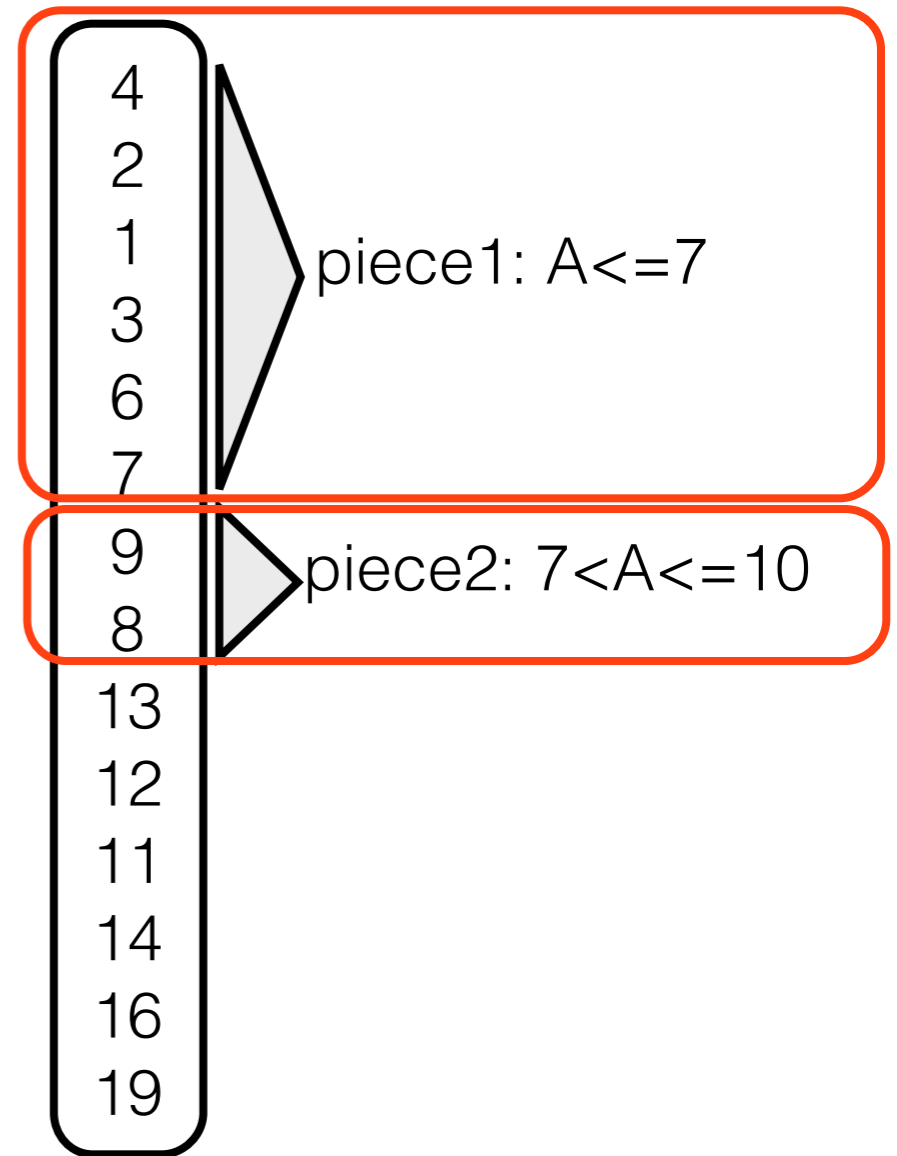
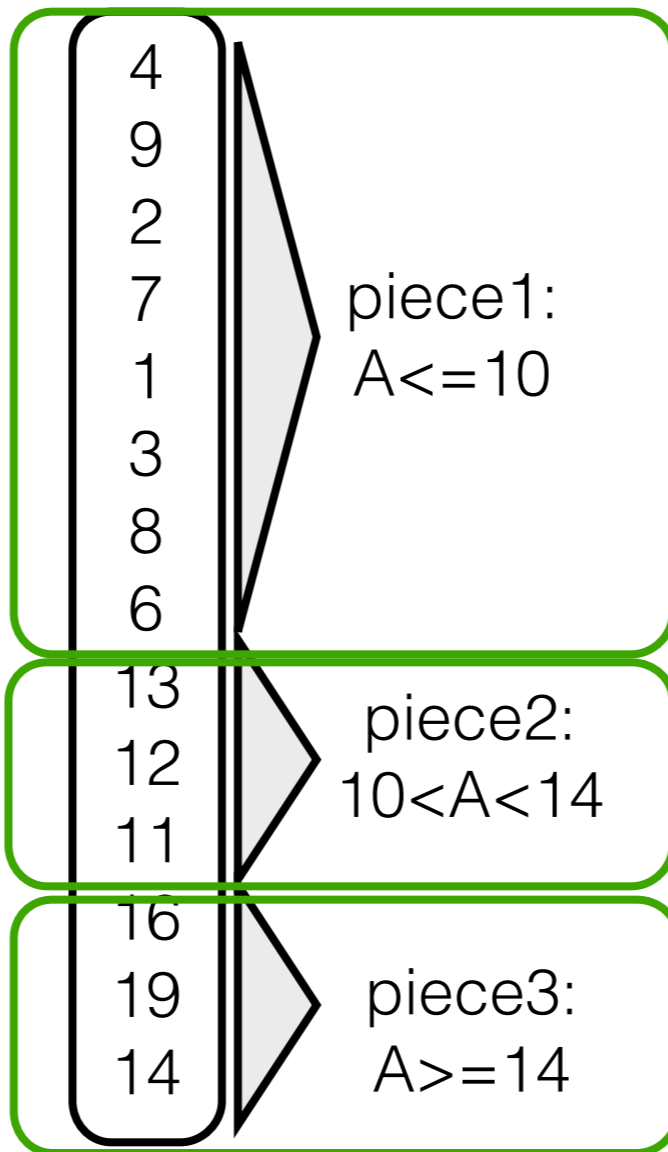
Database Cracking CIDR 2007

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



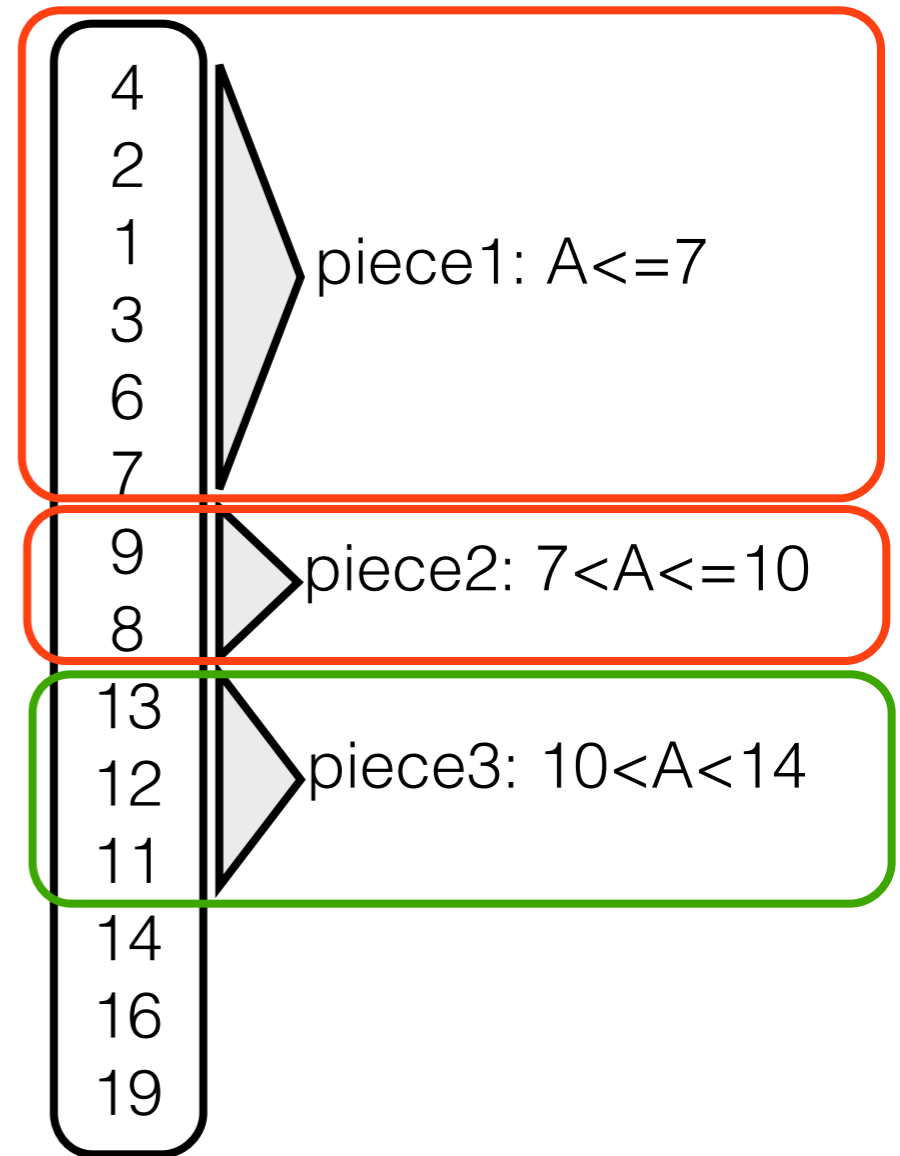
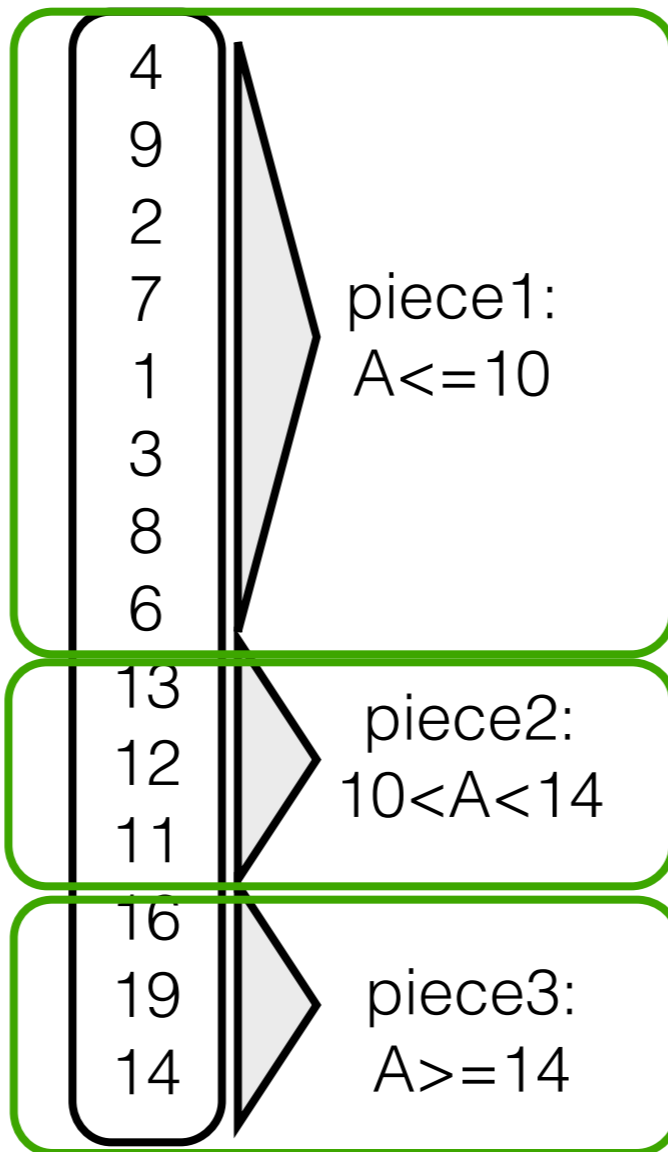
dynamically/on-the-fly within the select-operator

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

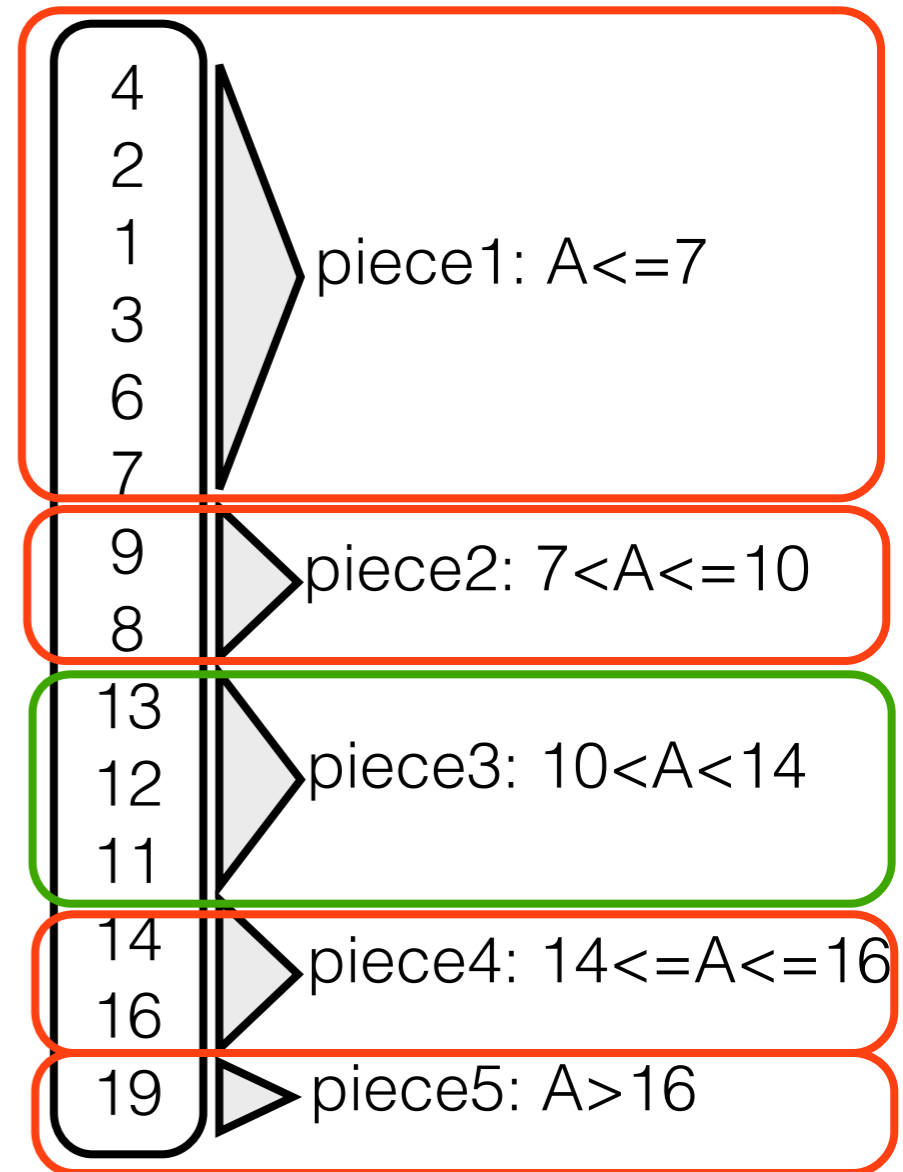
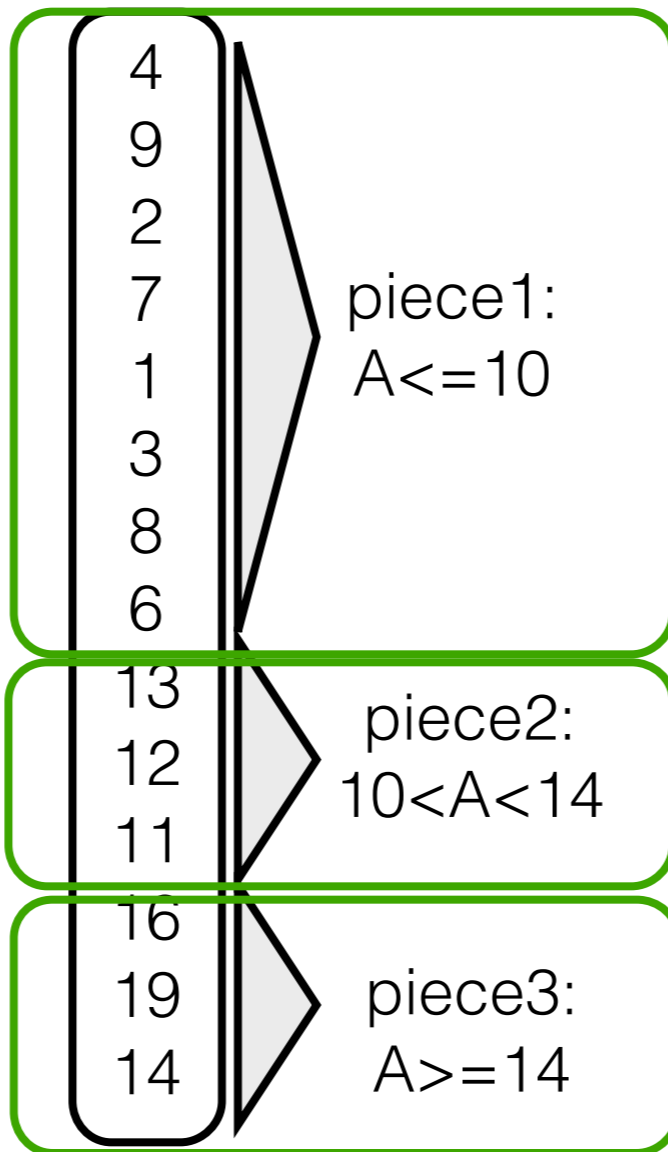
Database Cracking CIDR 2007

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

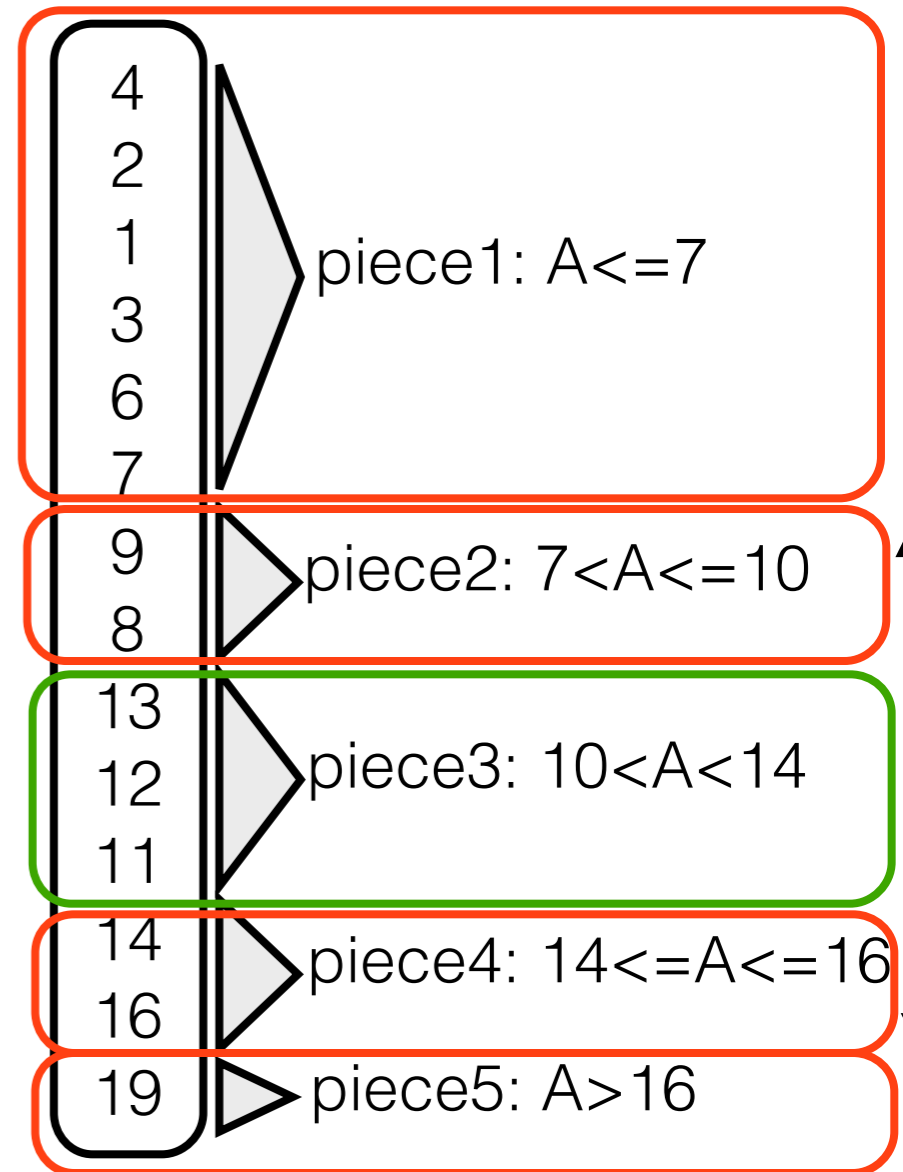
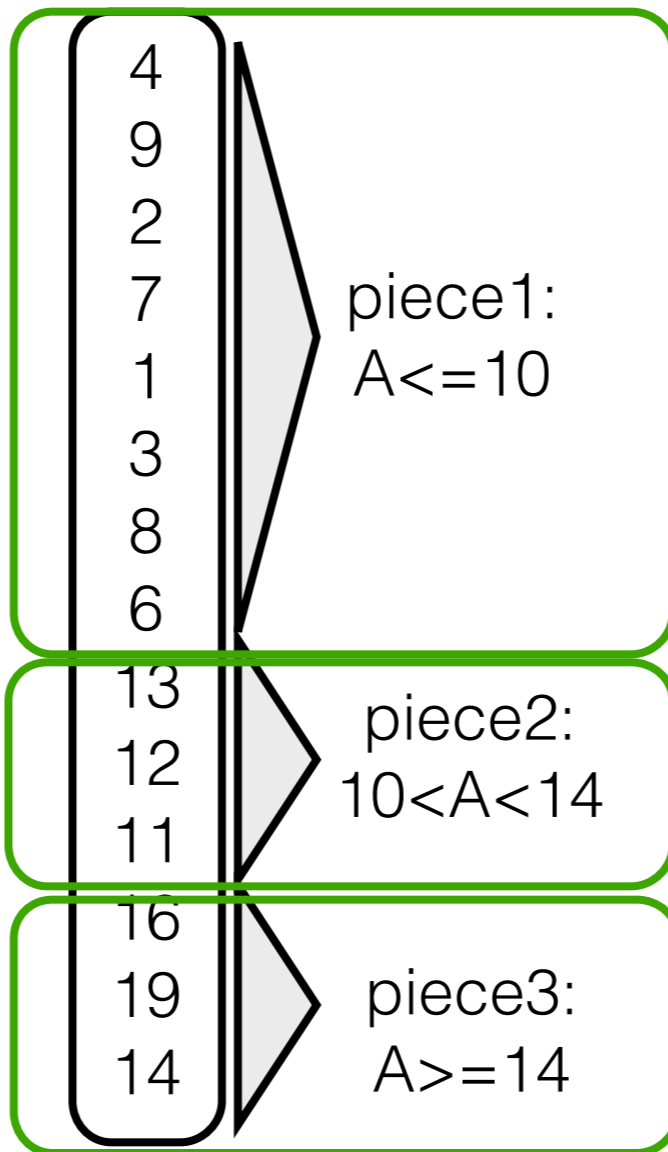
Database Cracking CIDR 2007

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

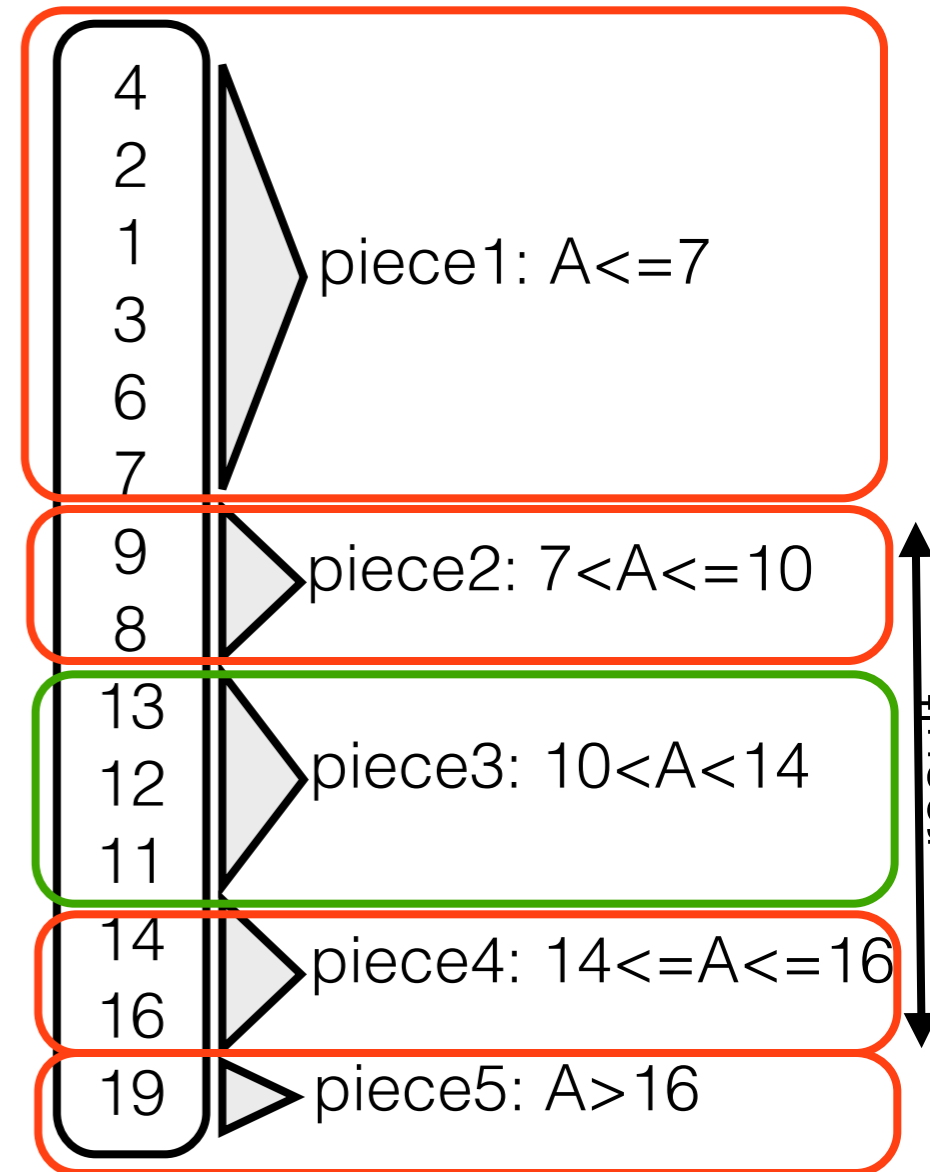
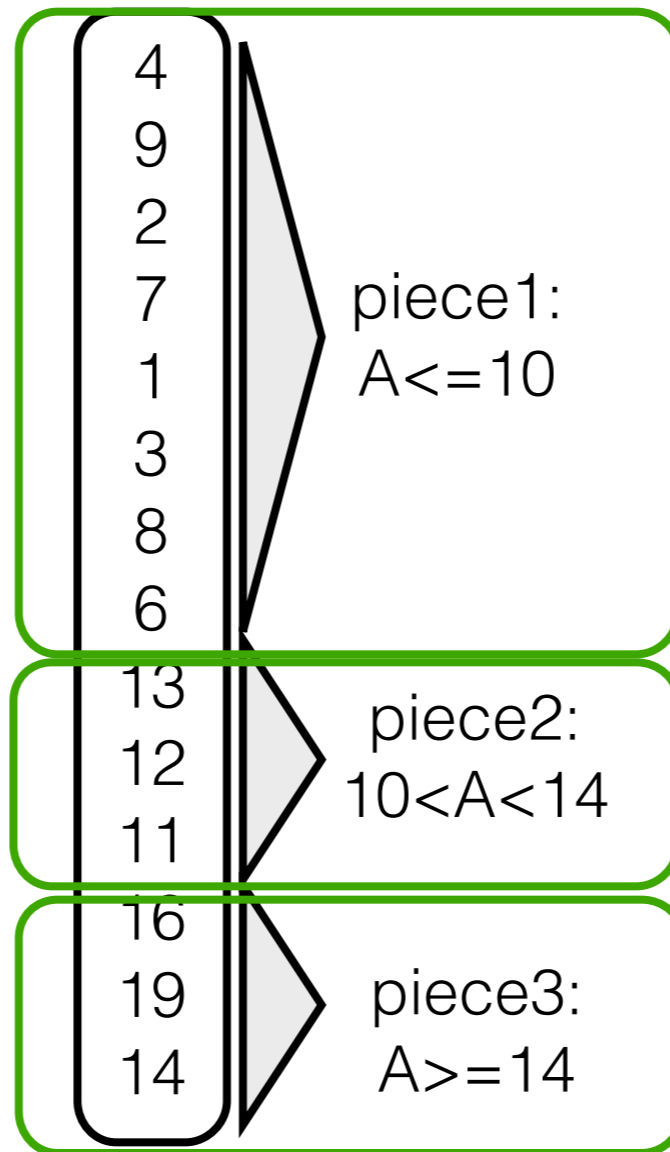
the more we crack, the more we learn

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

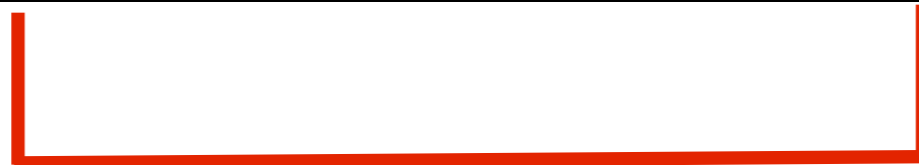


result

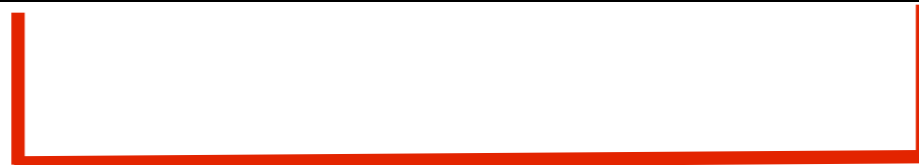
dynamically/on-the-fly within the select-operator

Database Cracking CIDR 2007

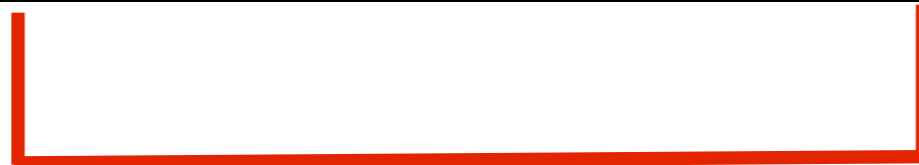




select [15,55]



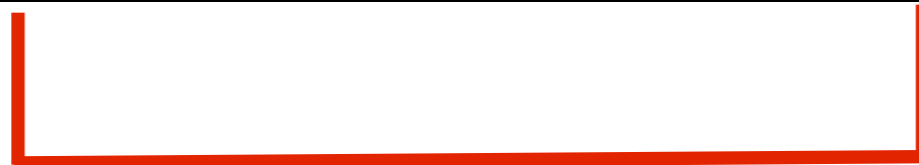
select [15,55]



select [15,55]

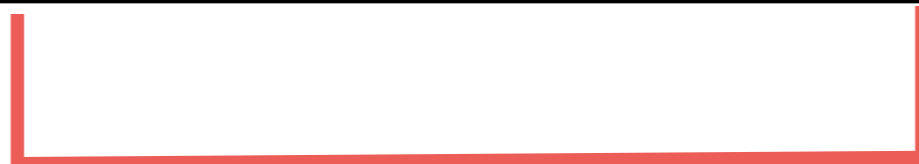
10 20 30 40 50 60



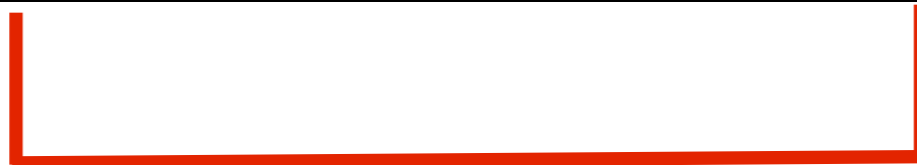


select [15,55]

10 20 30 40 50 60

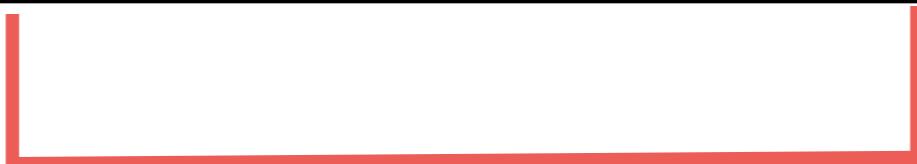


select [15,55]



select [15,55]

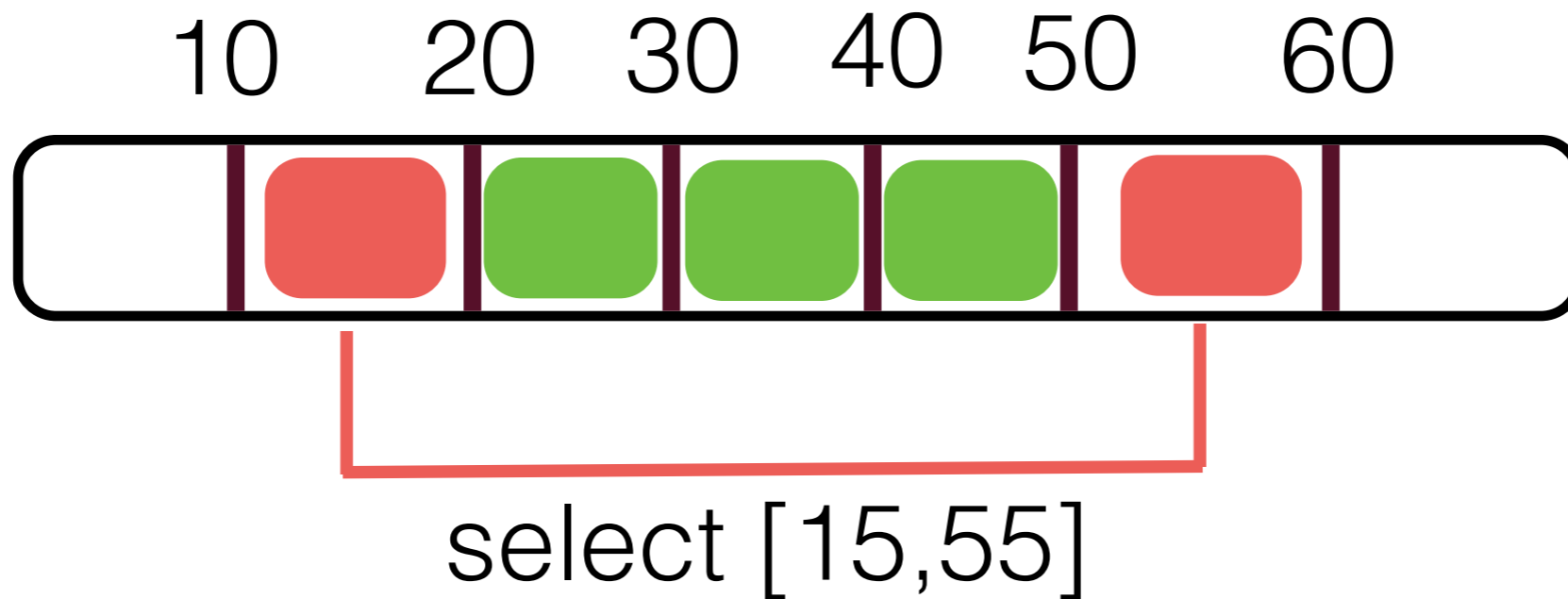
10 20 30 40 50 60



select [15,55]

touch at most two pieces at a time

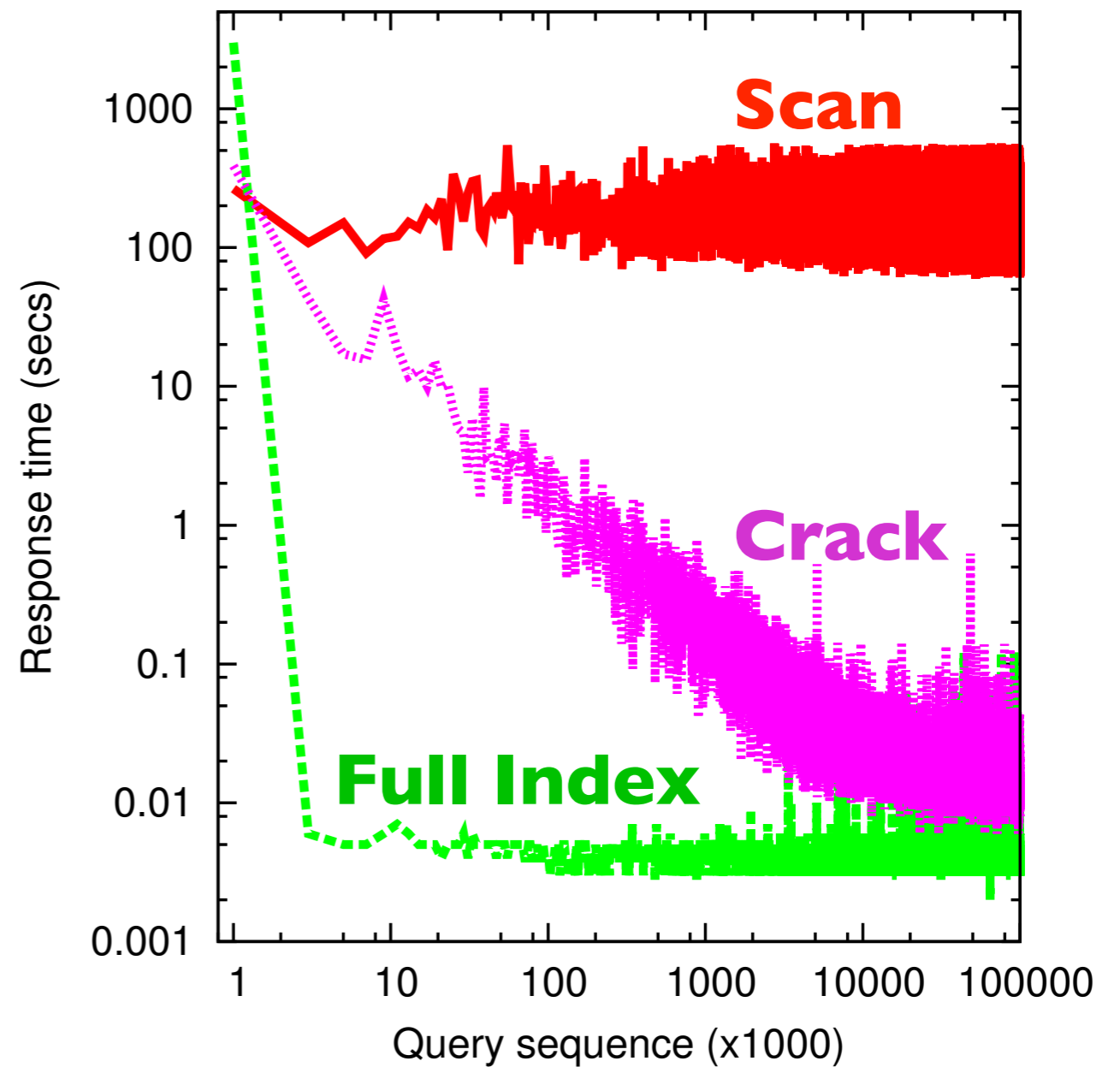
pieces become smaller and smaller



continuous adaptation

set-up

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

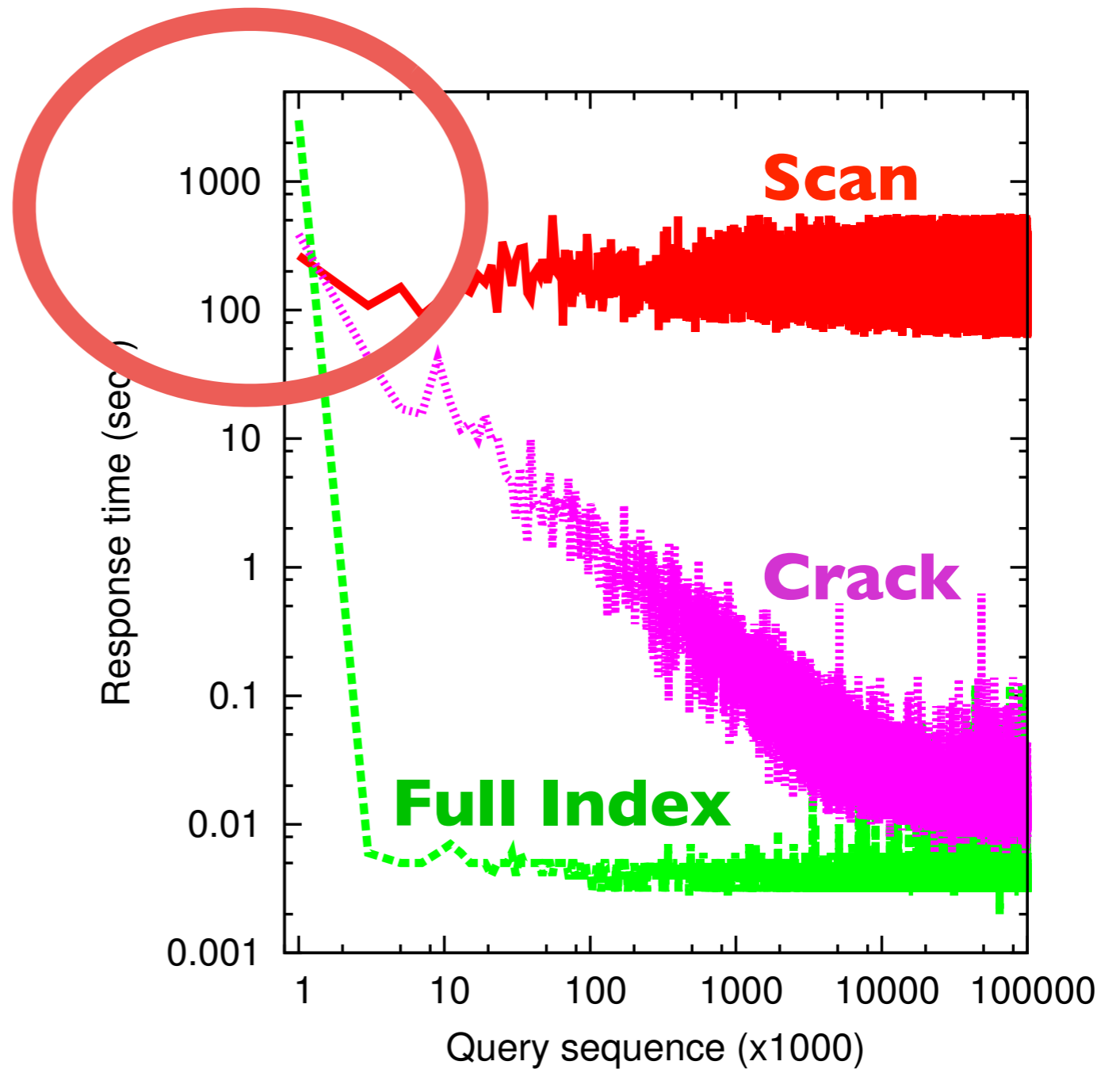


continuous adaptation

set-up

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

almost no
initialization overhead

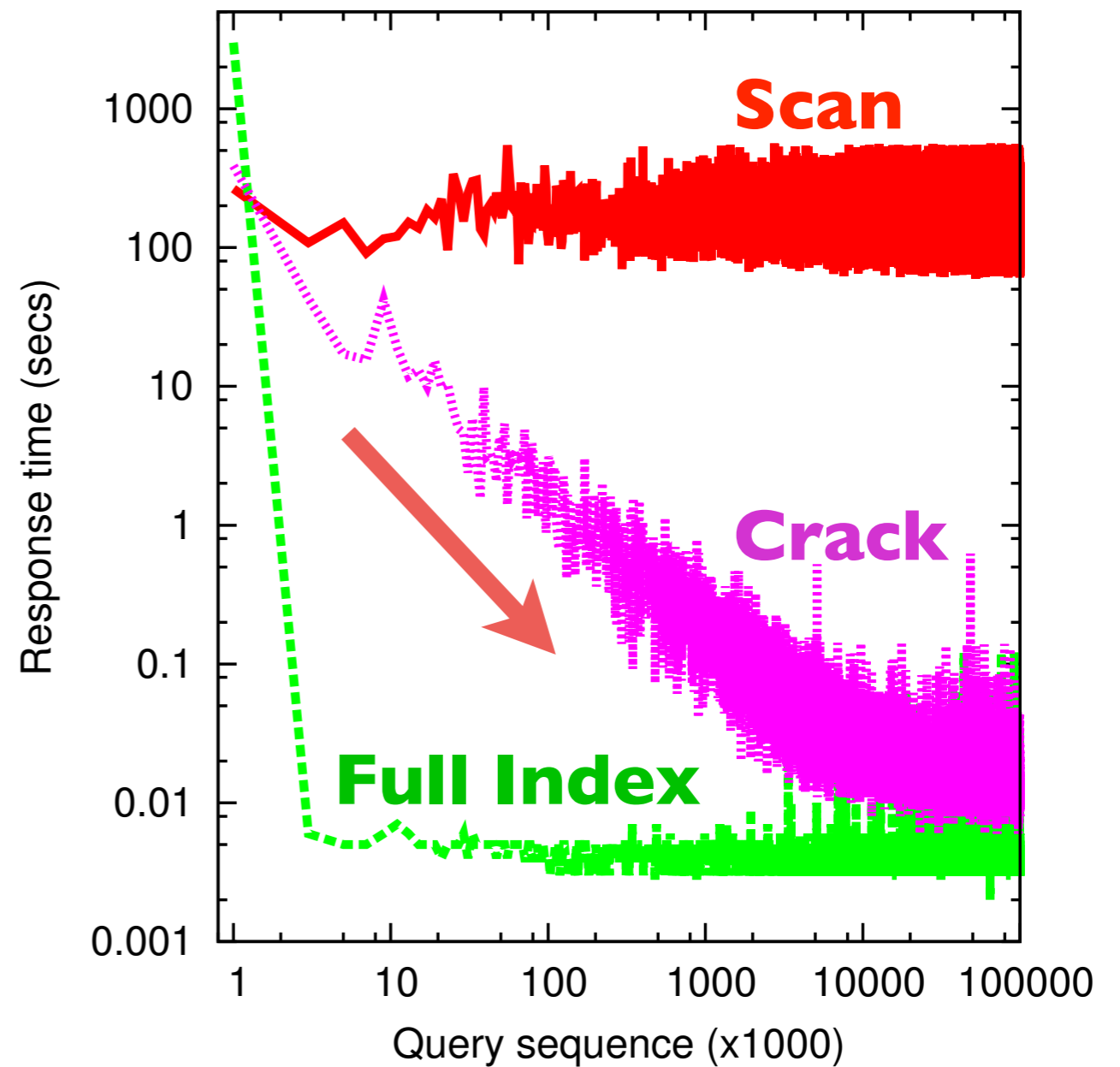


continuous adaptation

set-up

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

almost no
initialization overhead
continuous improvement

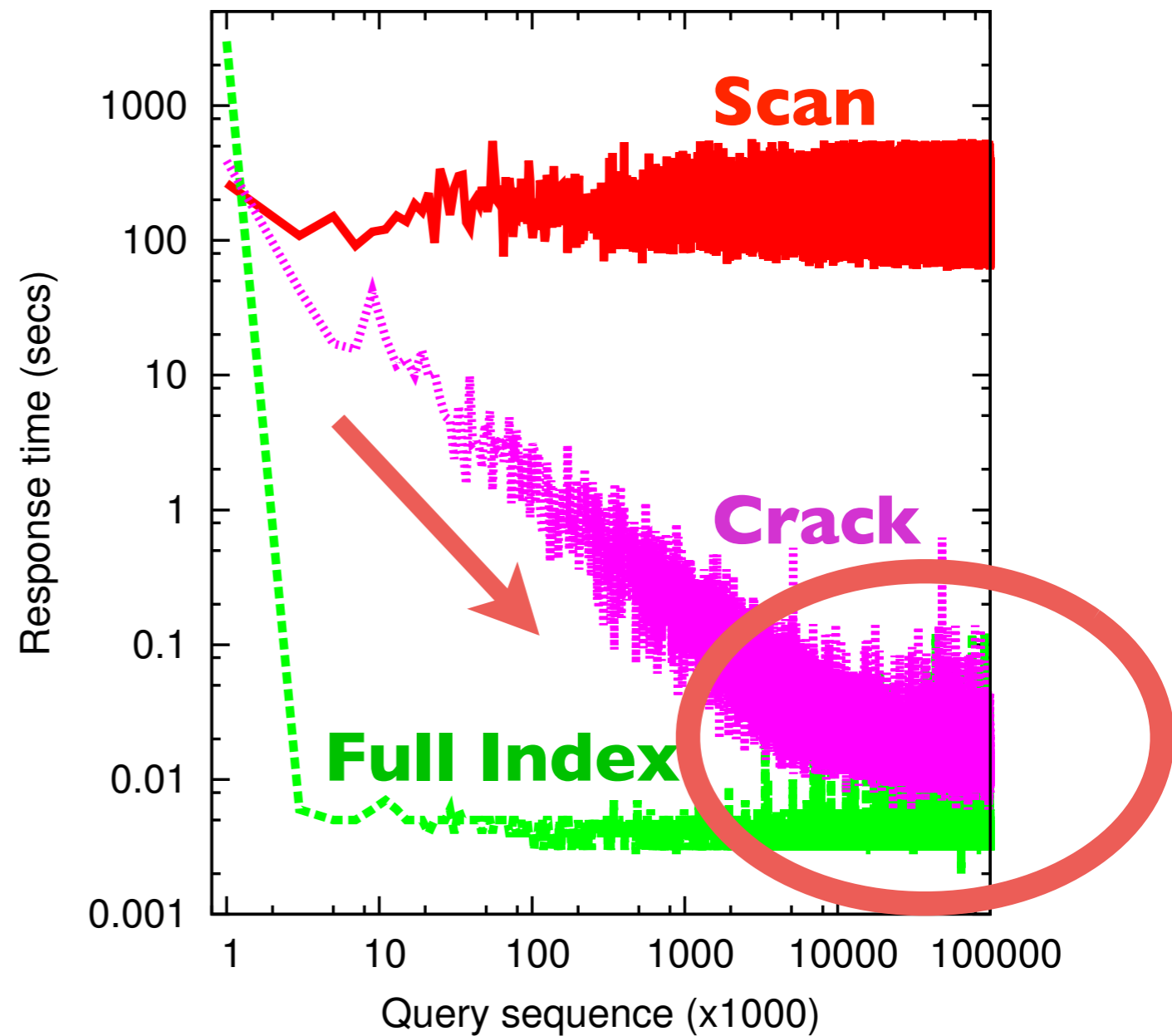


continuous adaptation

set-up

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

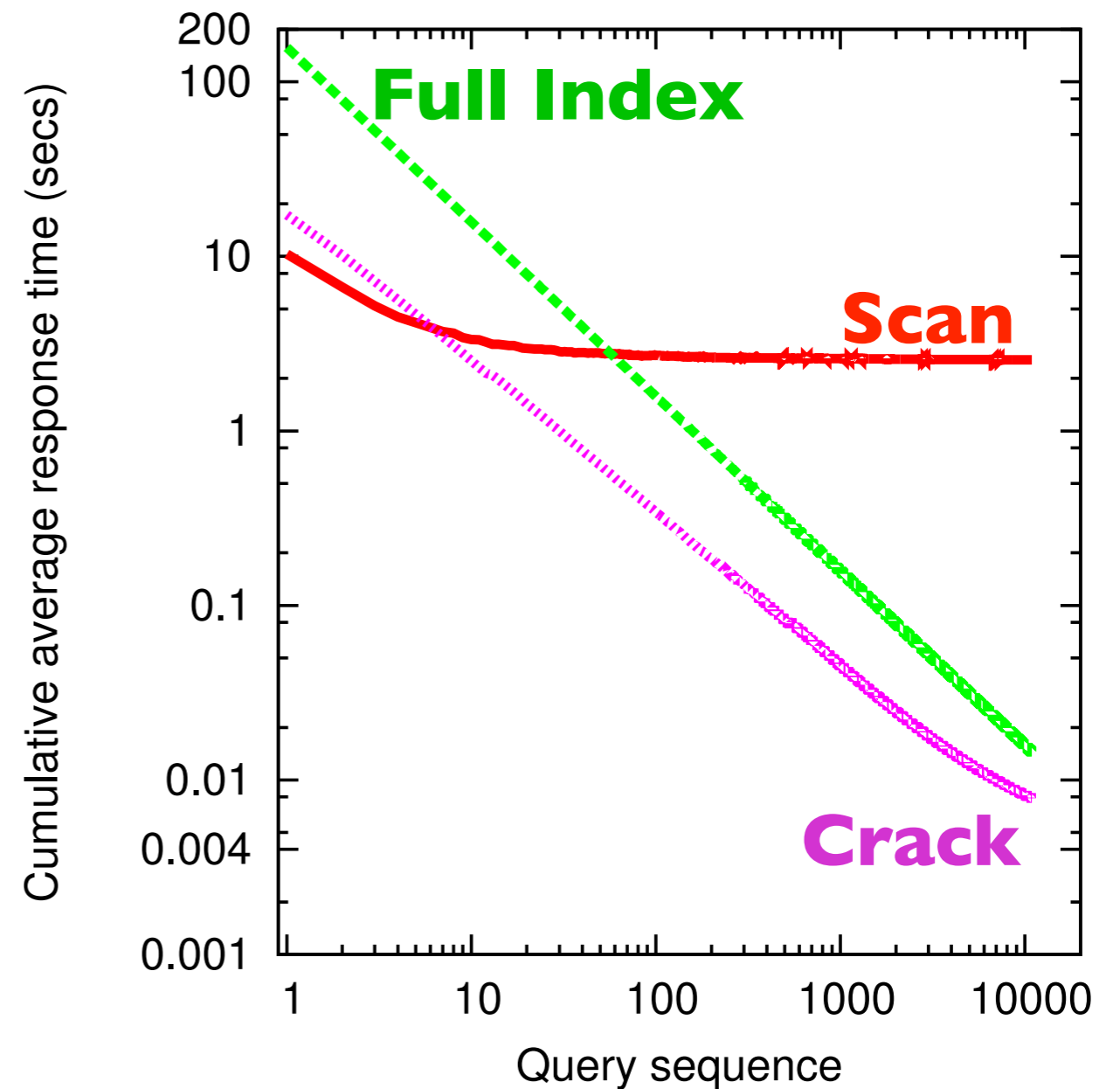
almost no
initialization overhead
continuous improvement



continuous adaptation

set-up

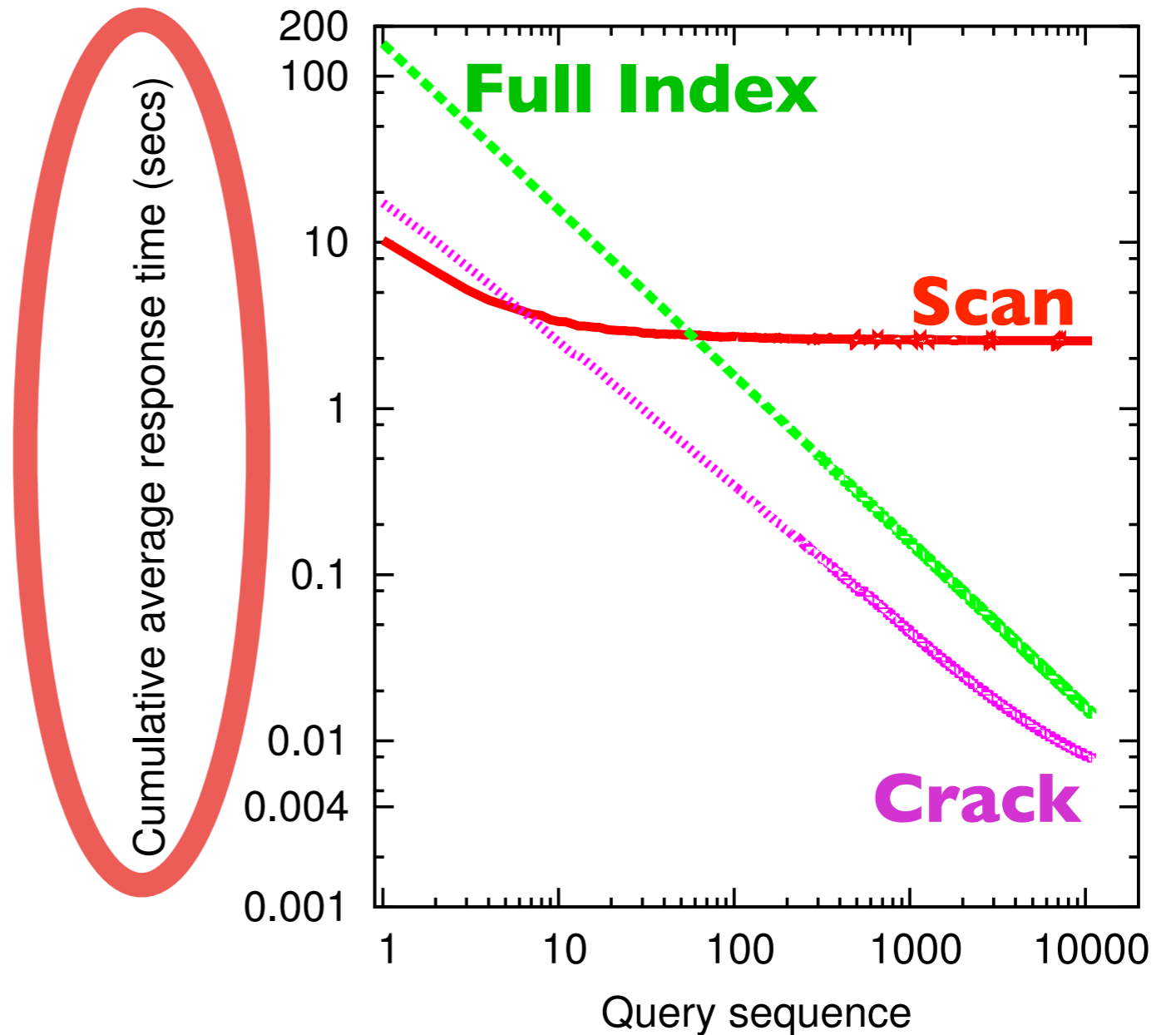
10K random selections
selectivity 10%
random value ranges
in a 30 million integer column



continuous adaptation

set-up

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



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

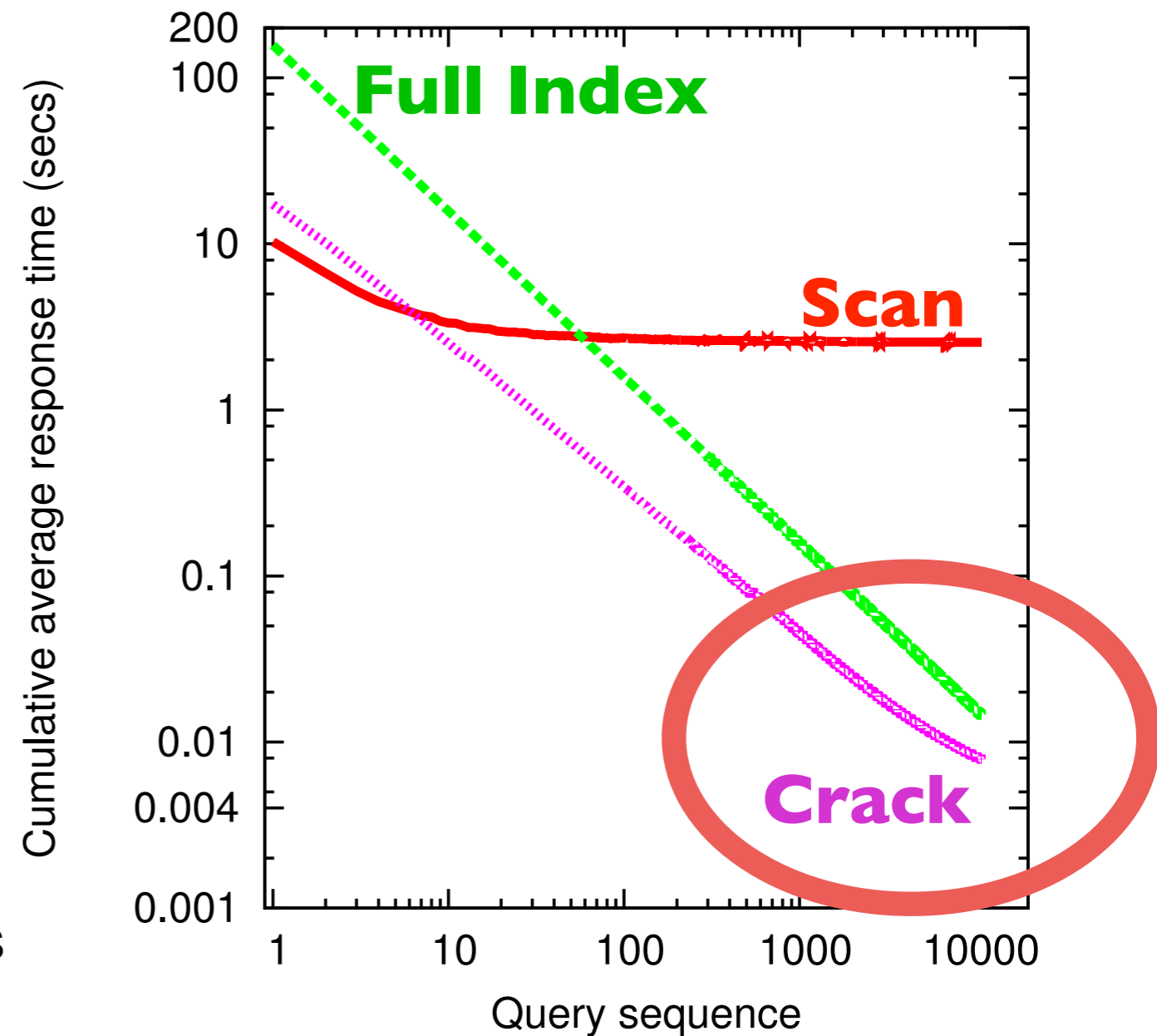


table 1

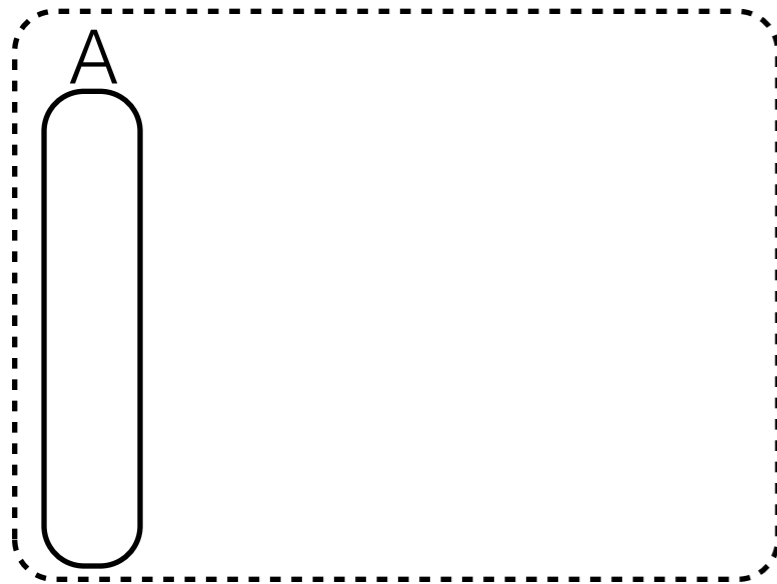
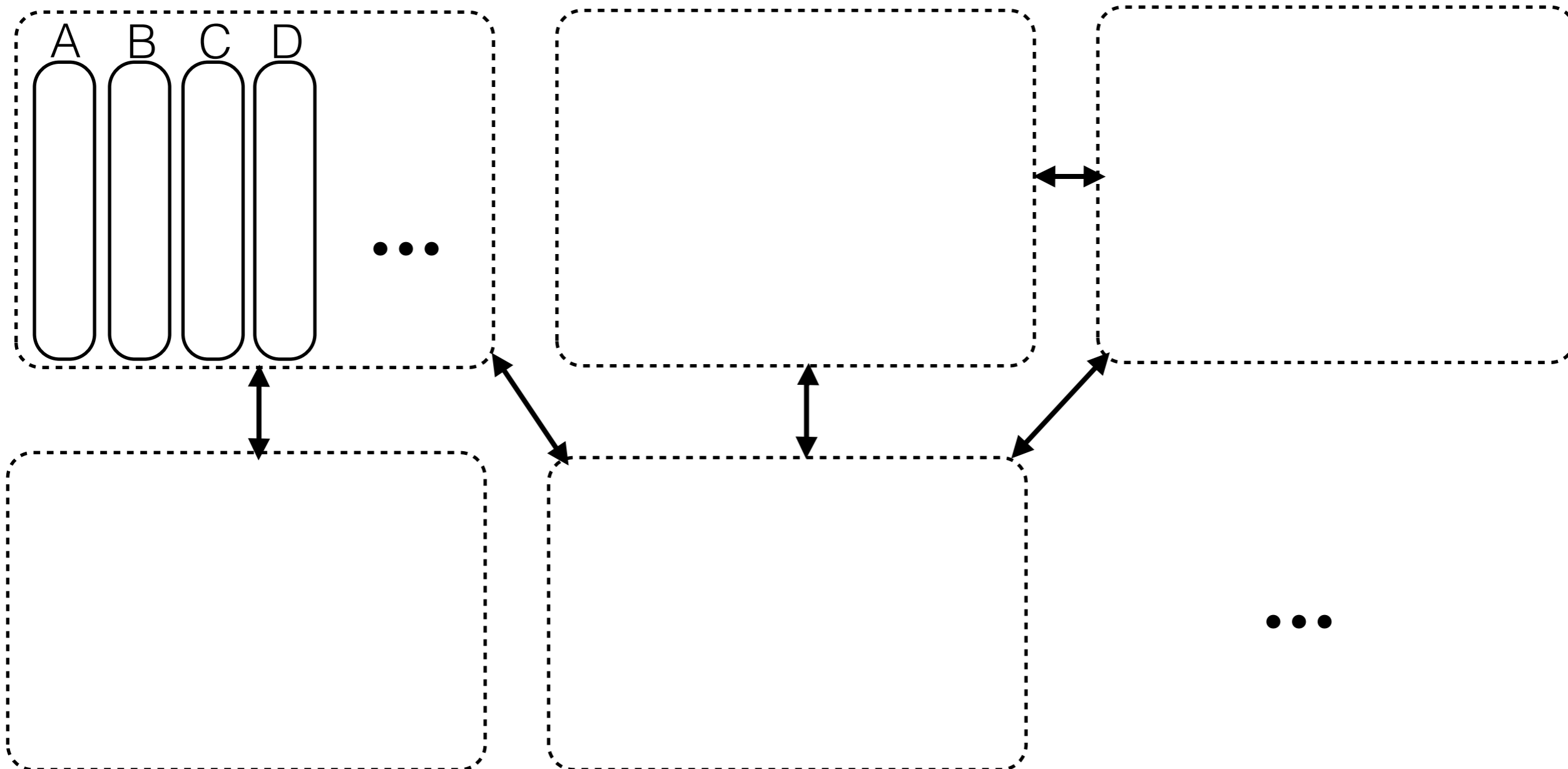
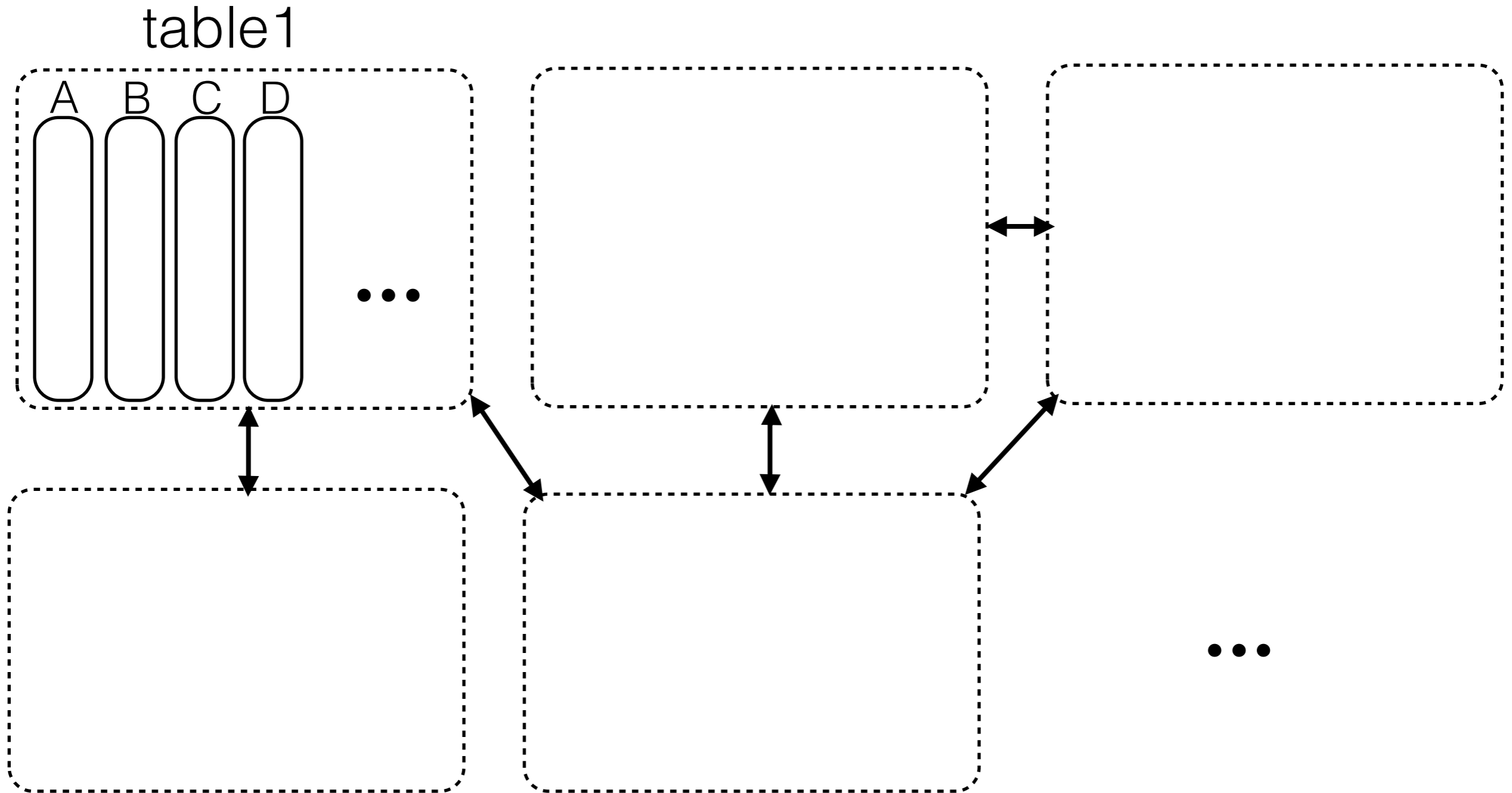


table 1



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

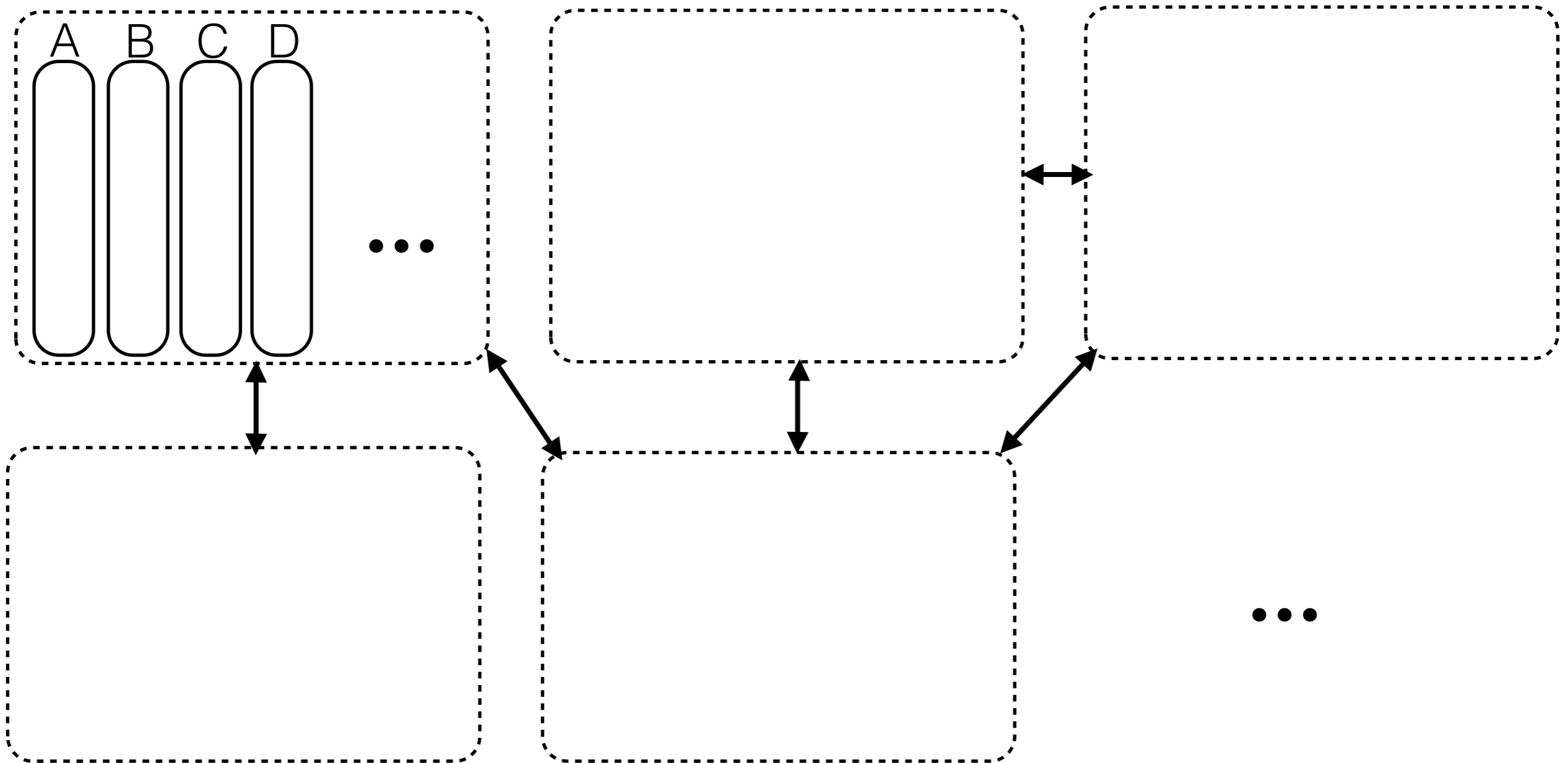
select max(R.A), max(R.B), max(S.A), max(S.B) from R, S

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



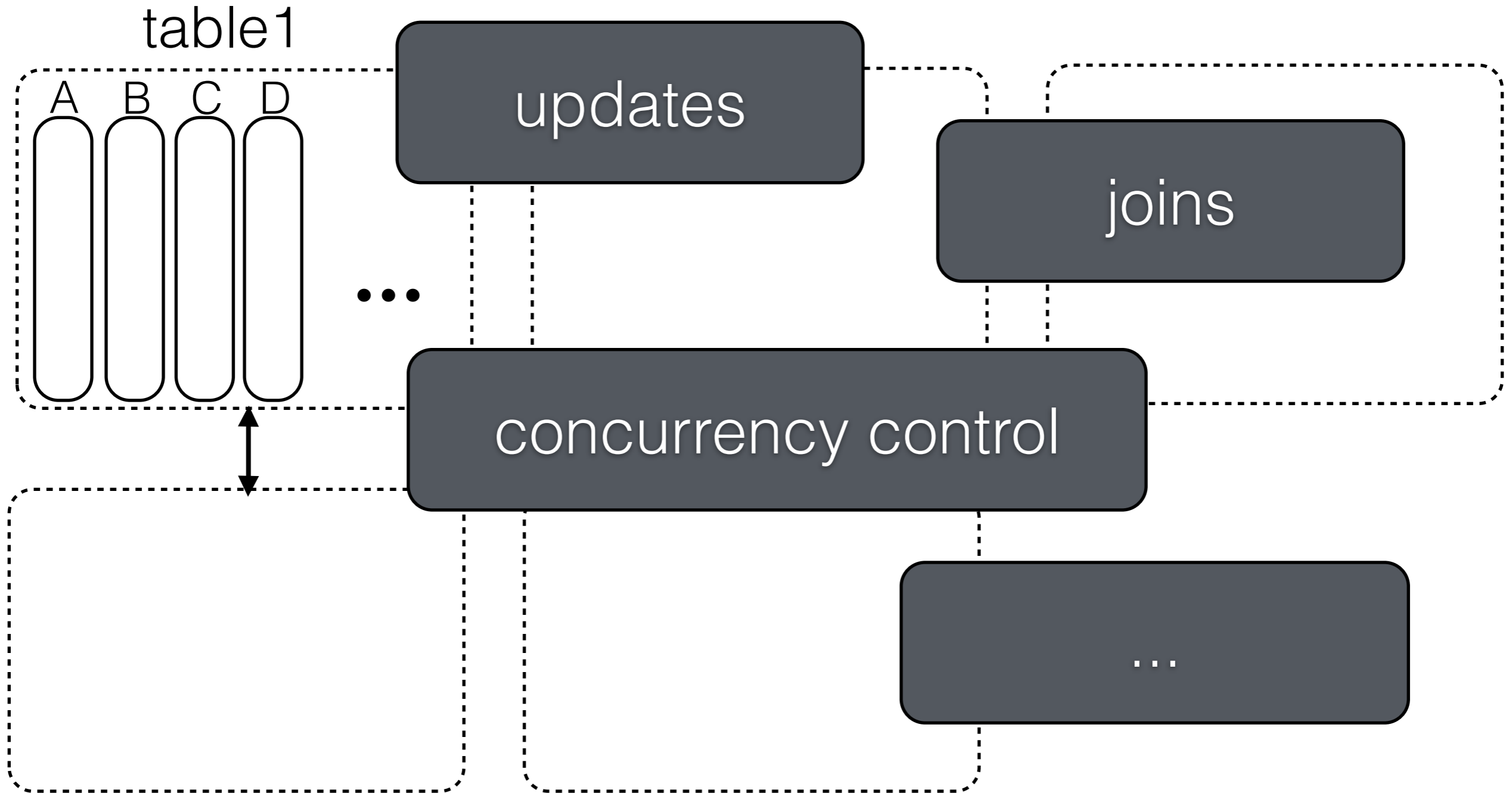
select R.A from R where R.A > 10 and R.A < 14

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

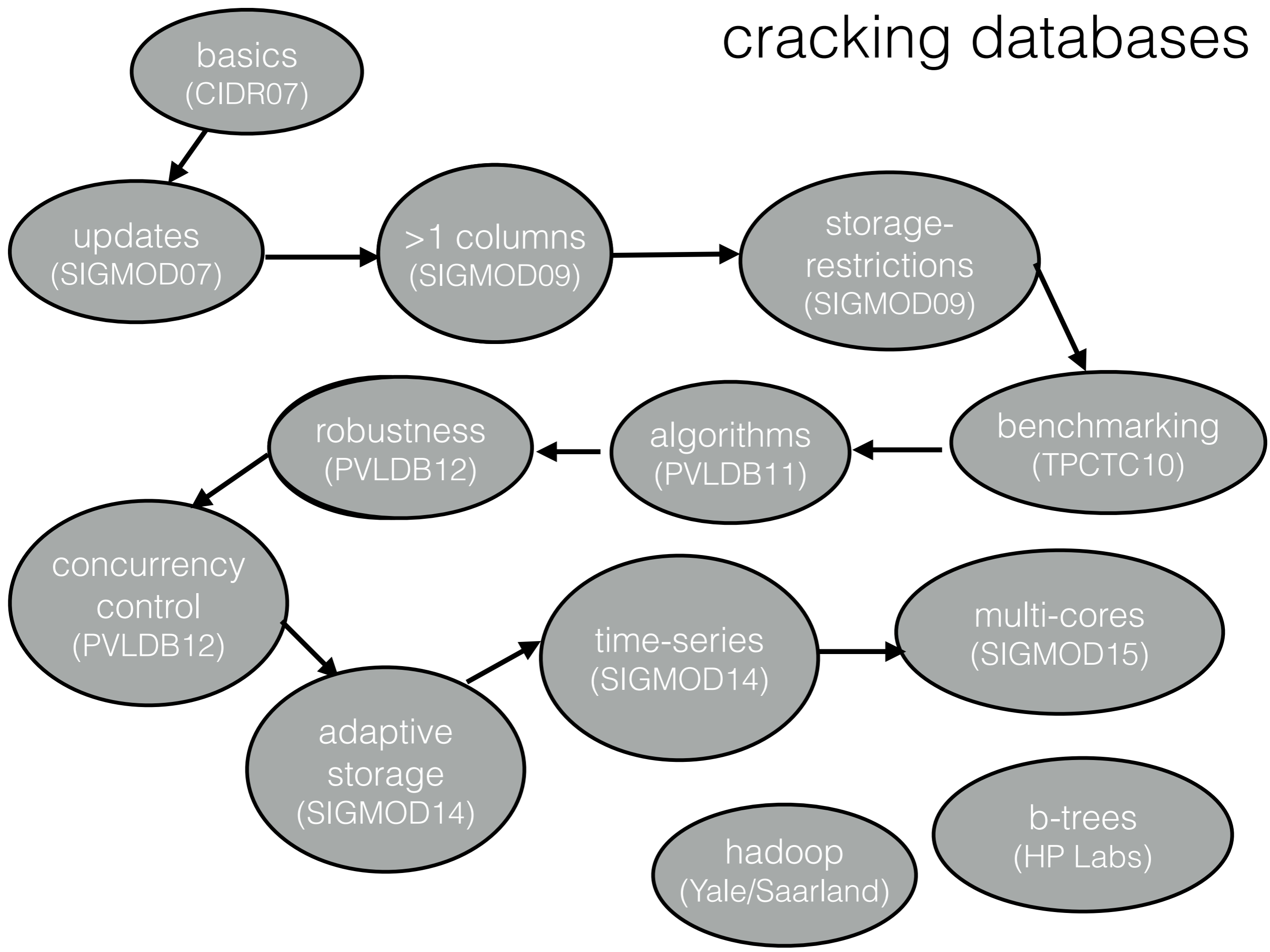
where v1 < R.C < v2 and v3 < R.D < v4

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

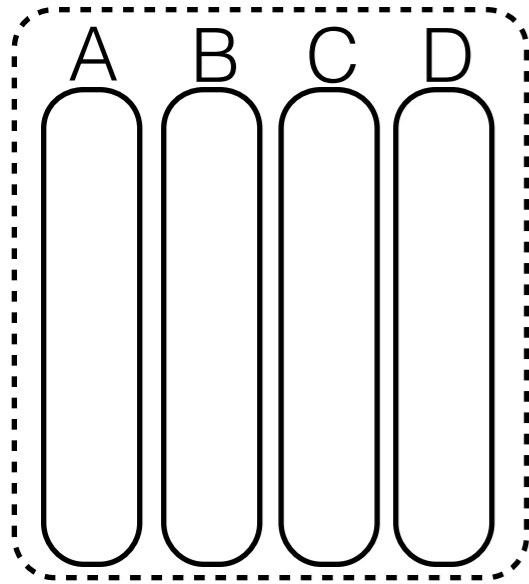


cracking databases



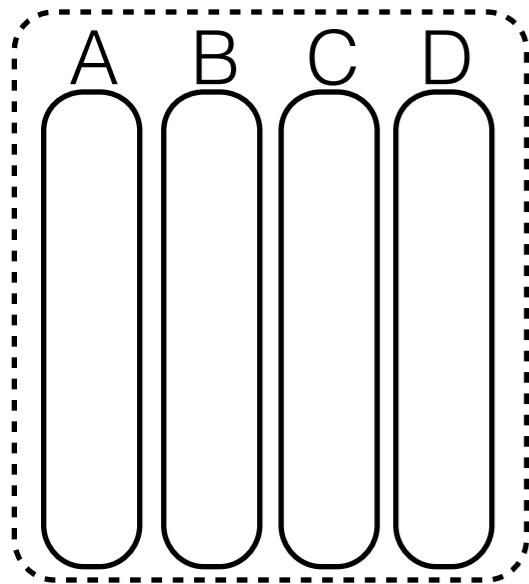
cracking tangram

base data
table 1



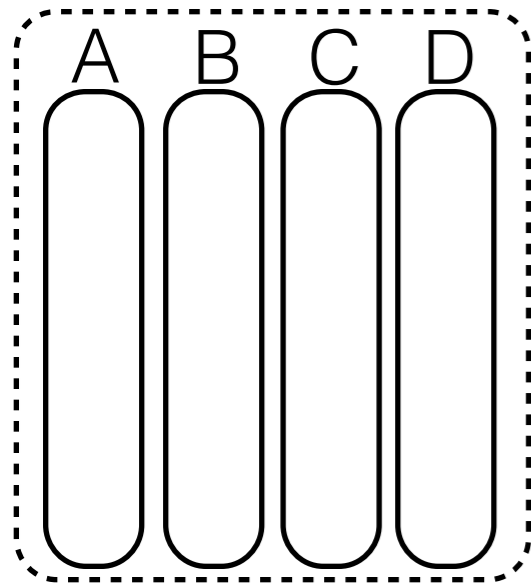
as queries arrive...

table 2



cracking tangram

base data
table 1



as queries arrive...
table 1

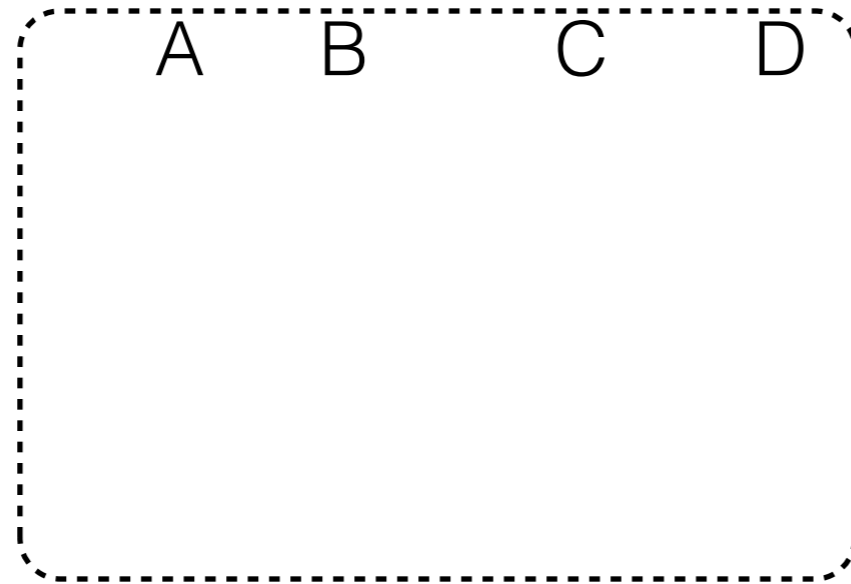


table 2

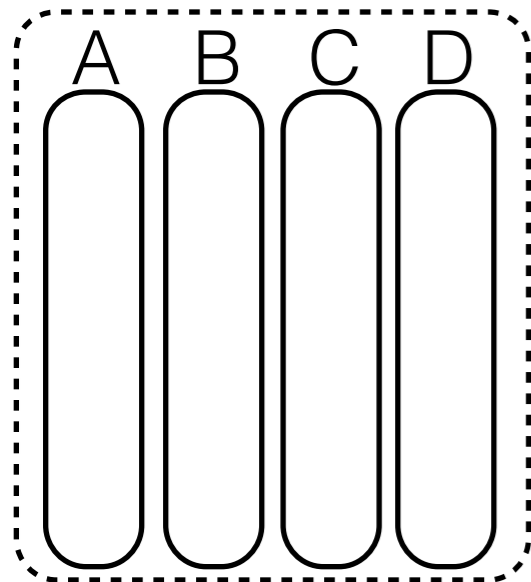
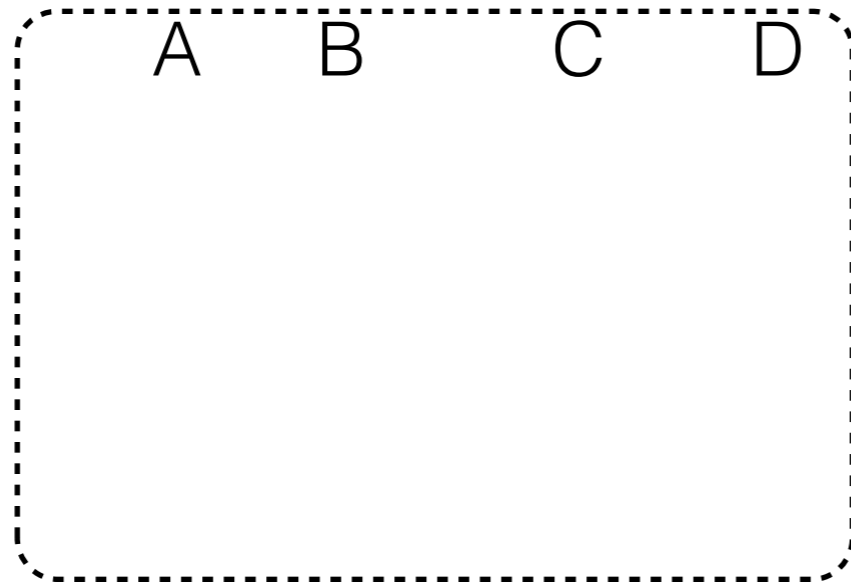
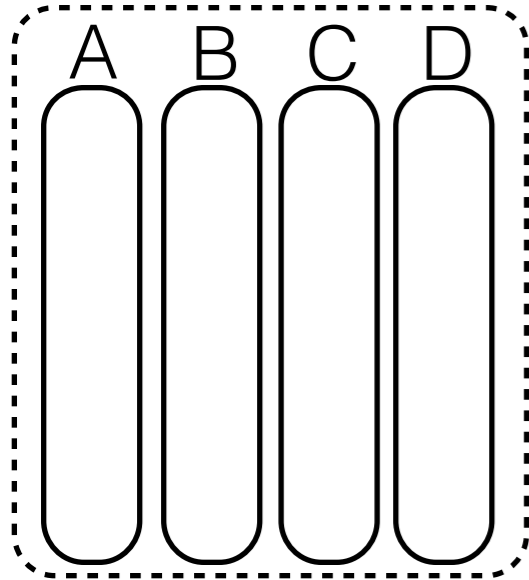


table 2

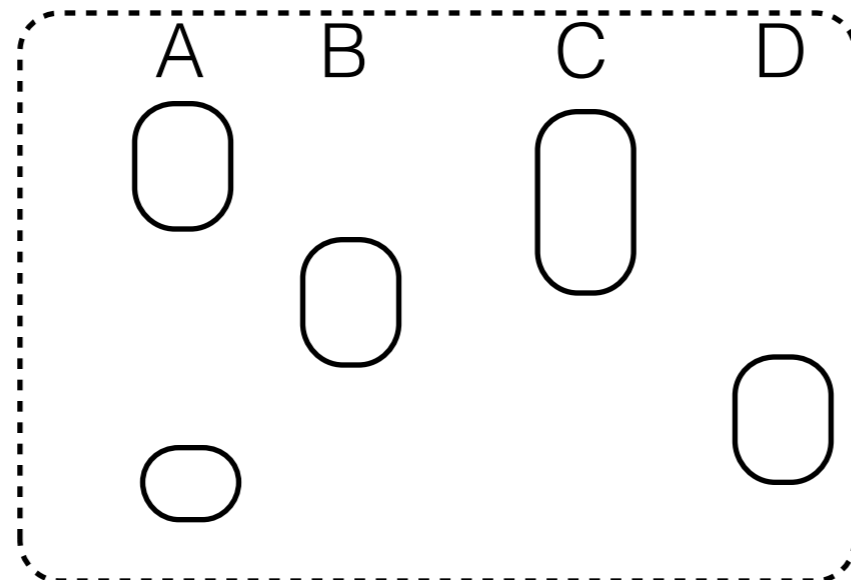


cracking tangram

base data
table 1



as queries arrive...
table 1



partial materialization

table 2

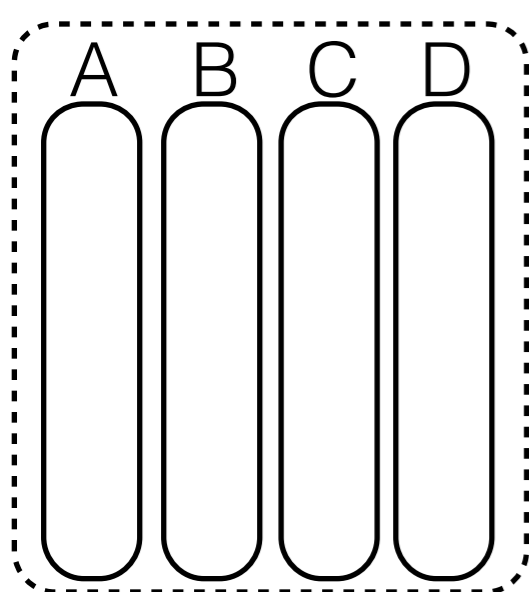
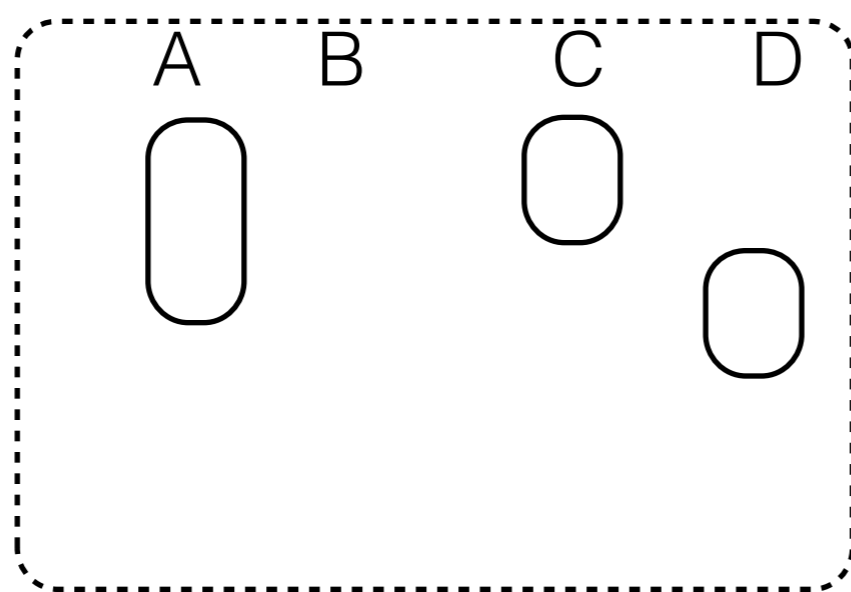
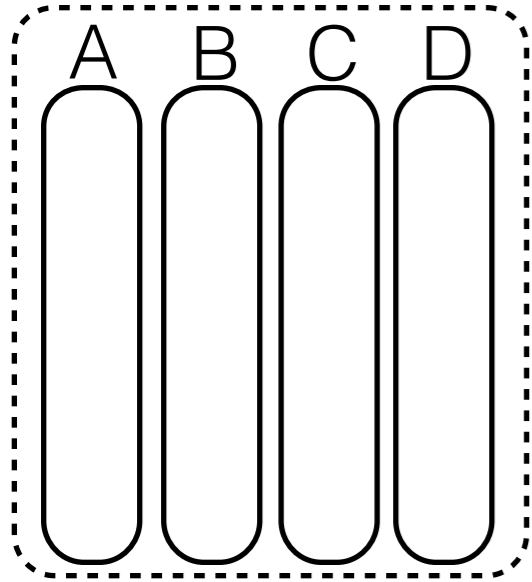


table 2

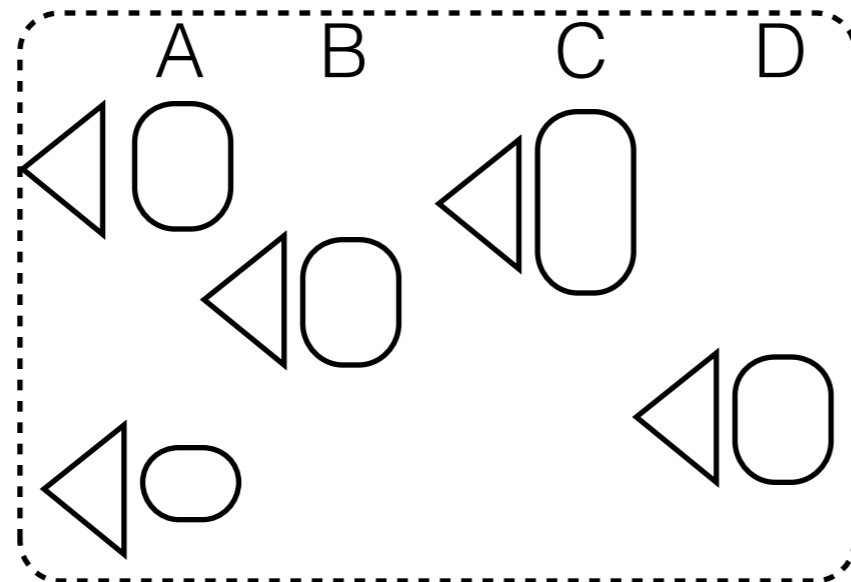


cracking tangram

base data
table 1



as queries arrive...
table 1



partial materialization
partial indexing

table 2

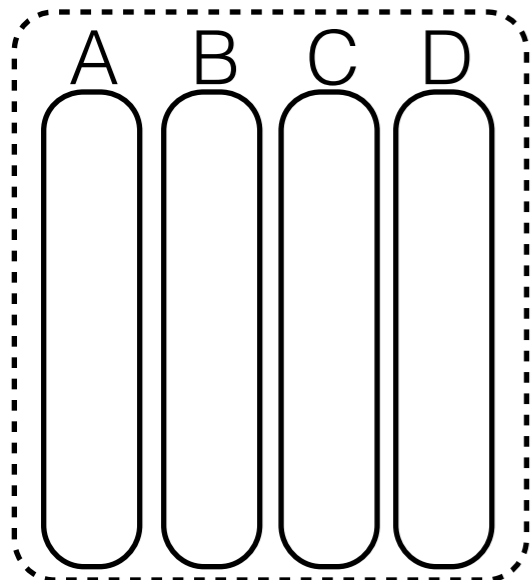
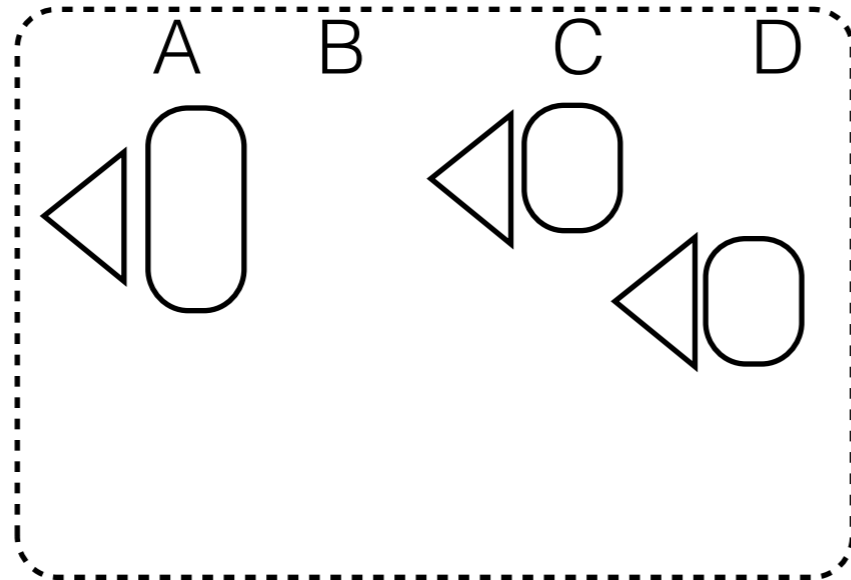
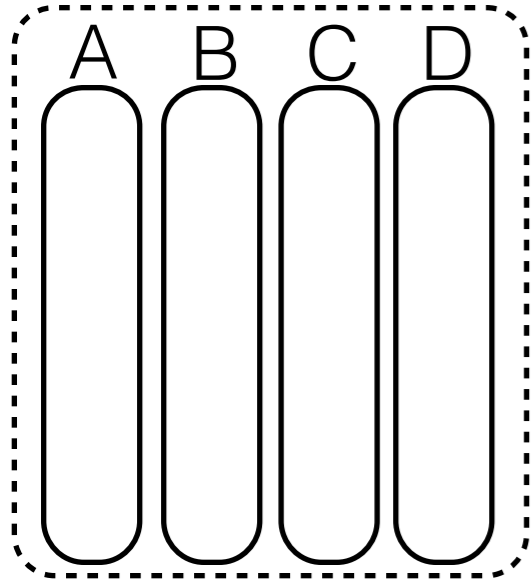


table 2



cracking tangram

base data
table 1



as queries arrive...
table 1

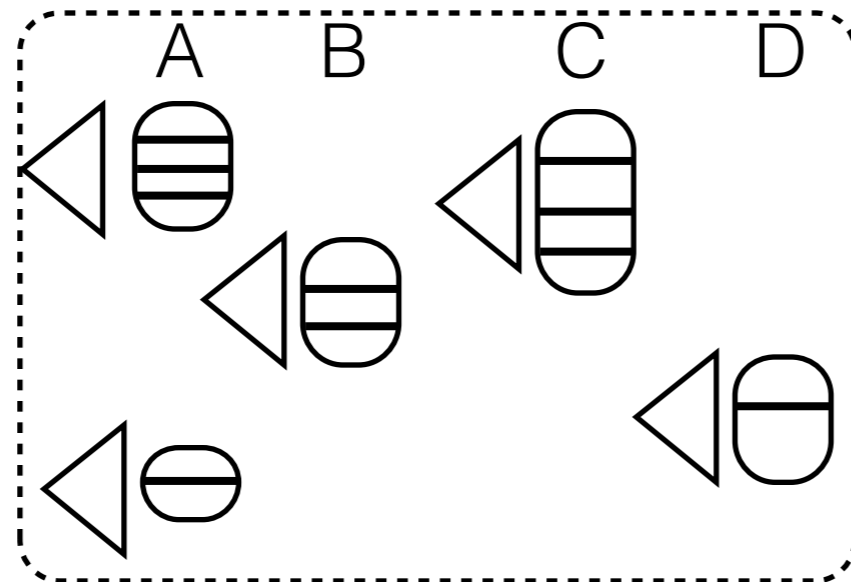


table 2

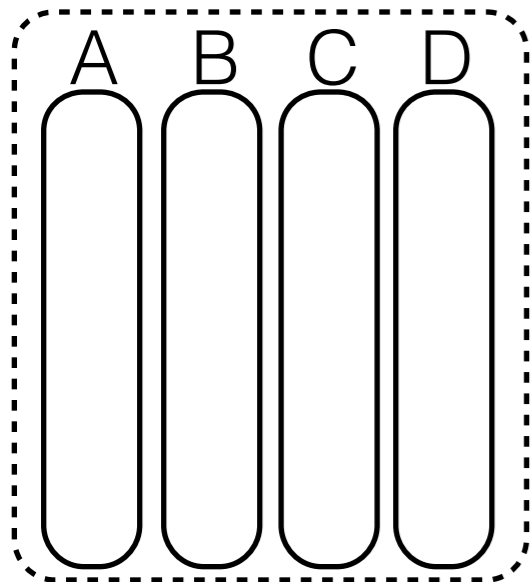
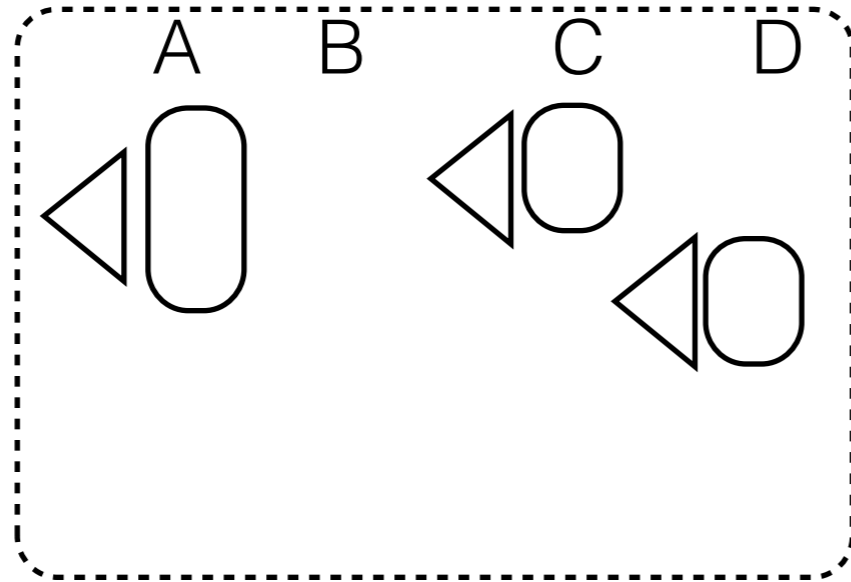


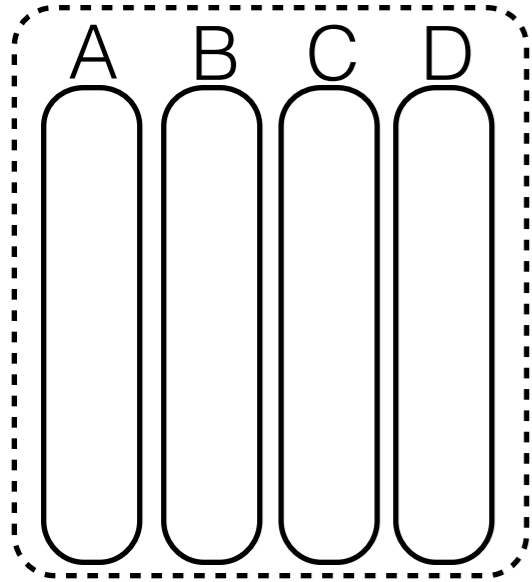
table 2



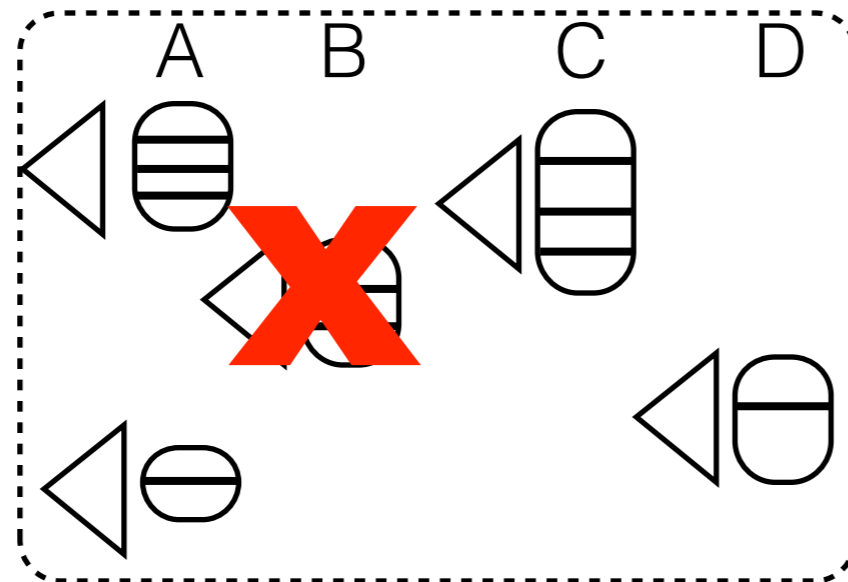
partial materialization
partial indexing
continuous adaptation

cracking tangram

base data
table 1



as queries arrive...
table 1



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

table 2

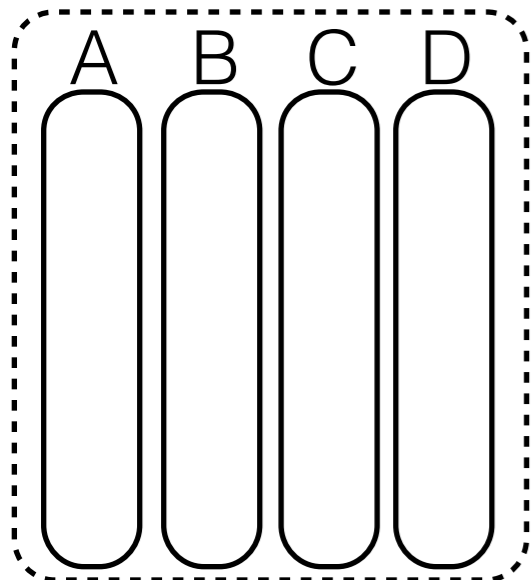
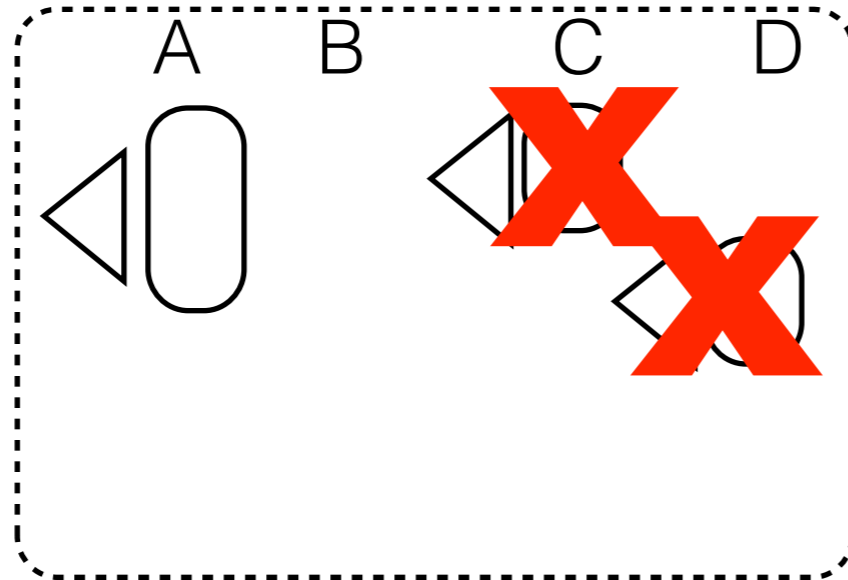
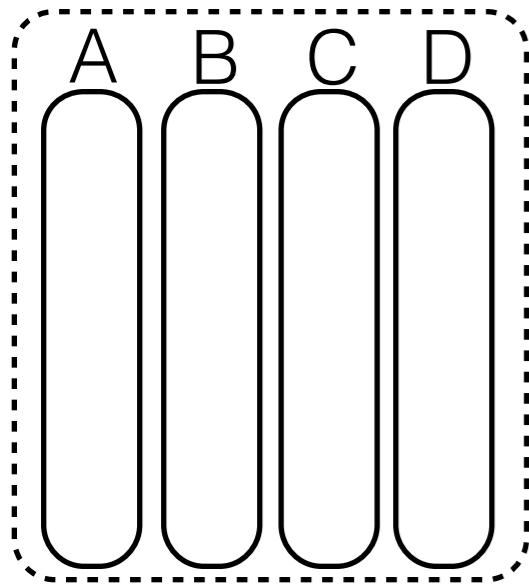


table 2

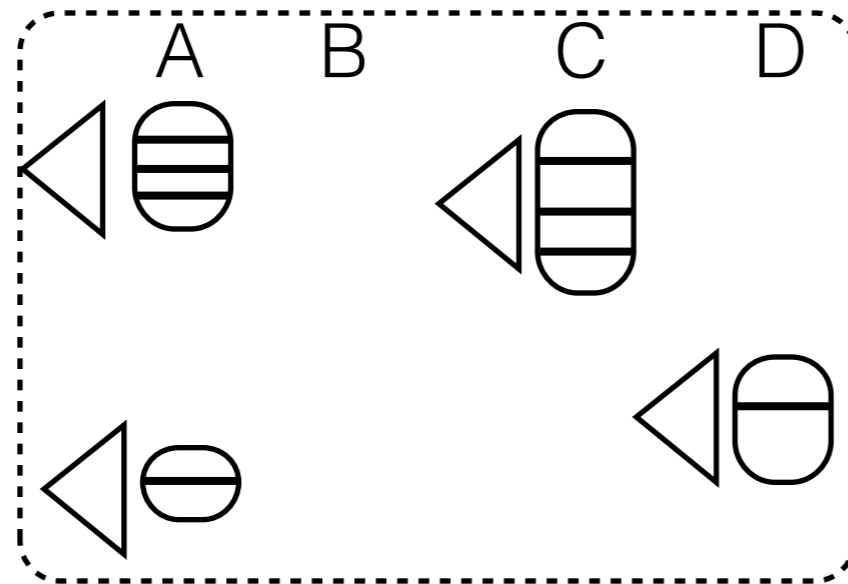


cracking tangram

base data
table 1



as queries arrive...
table 1



partial materialization
partial indexing
continuous adaptation
storage adaptation

table 2

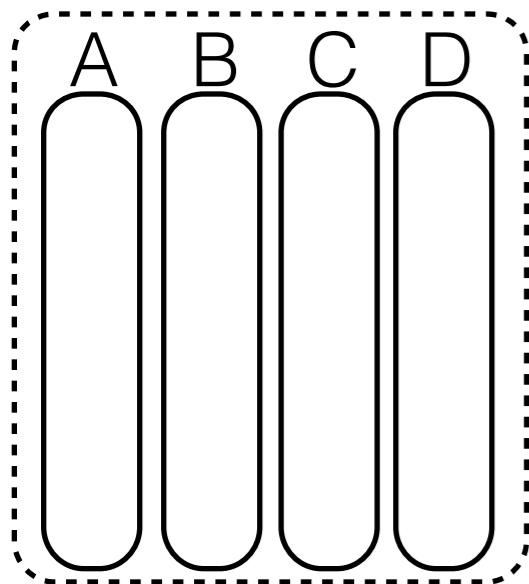
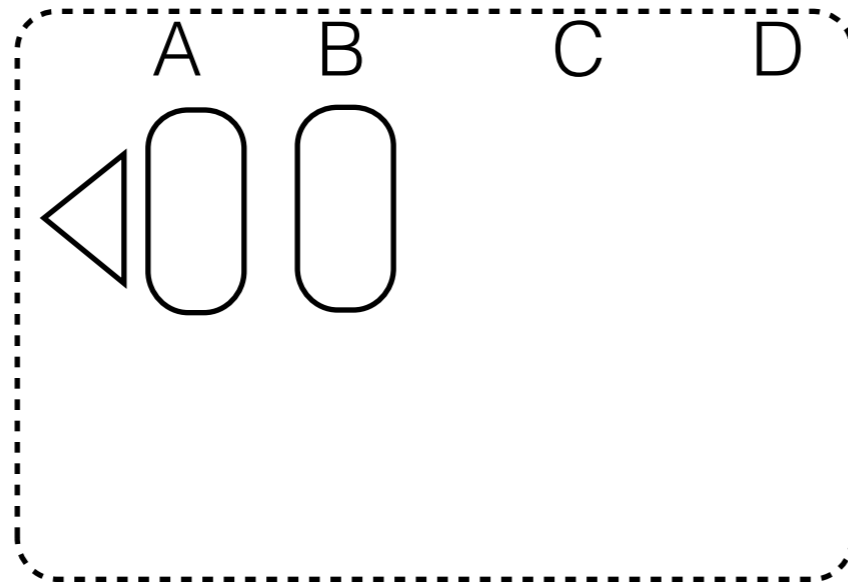
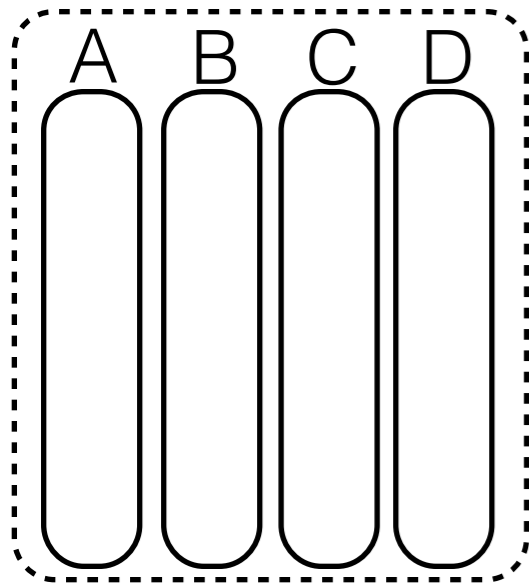


table 2



cracking tangram

base data
table 1



as queries arrive...
table 1

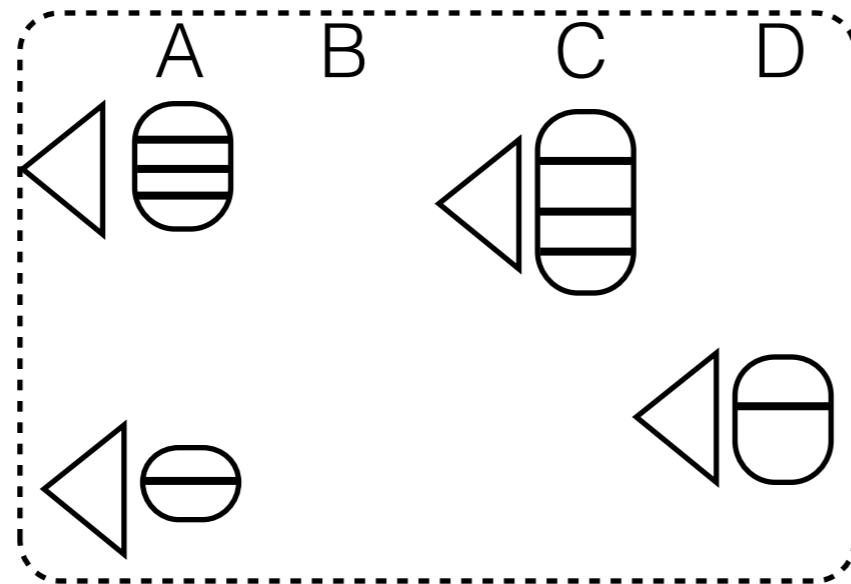


table 2

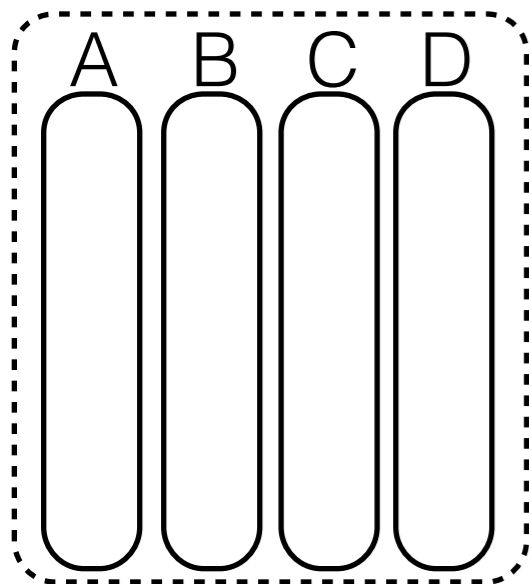
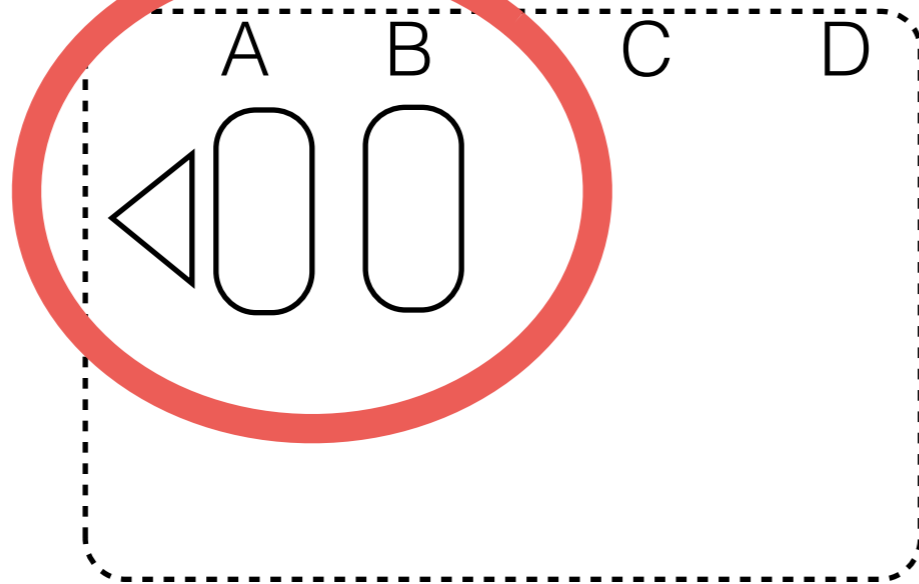


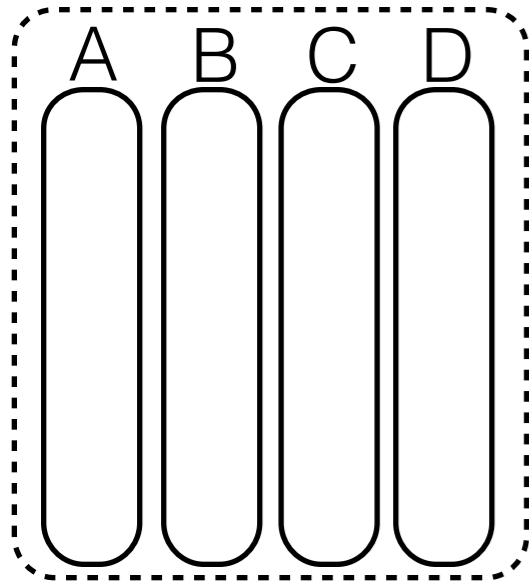
table 2



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

cracking tangram

base data
table 1



as queries arrive...
table 1

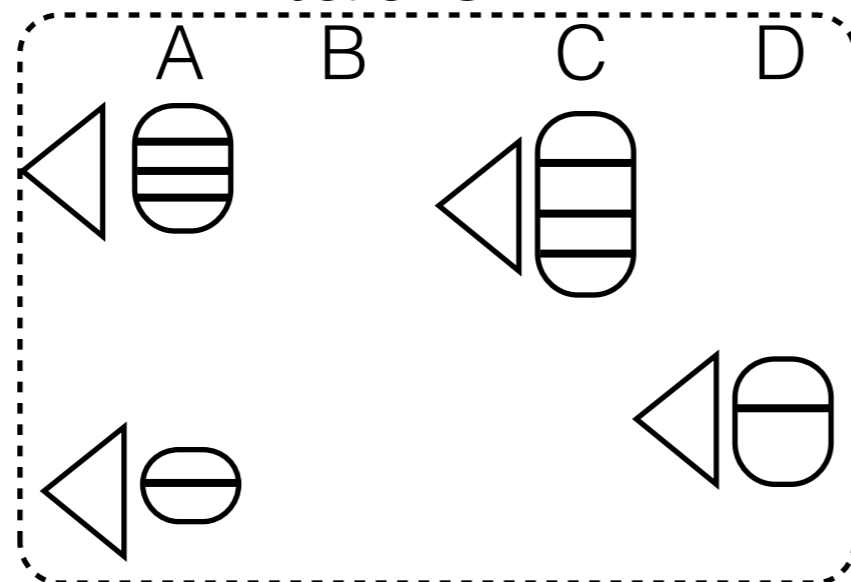


table 2

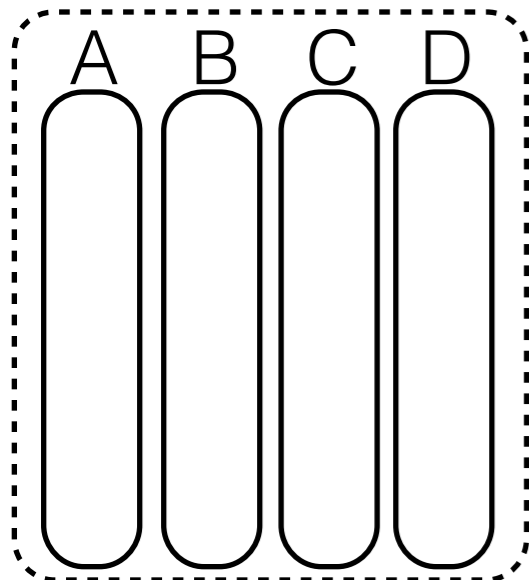
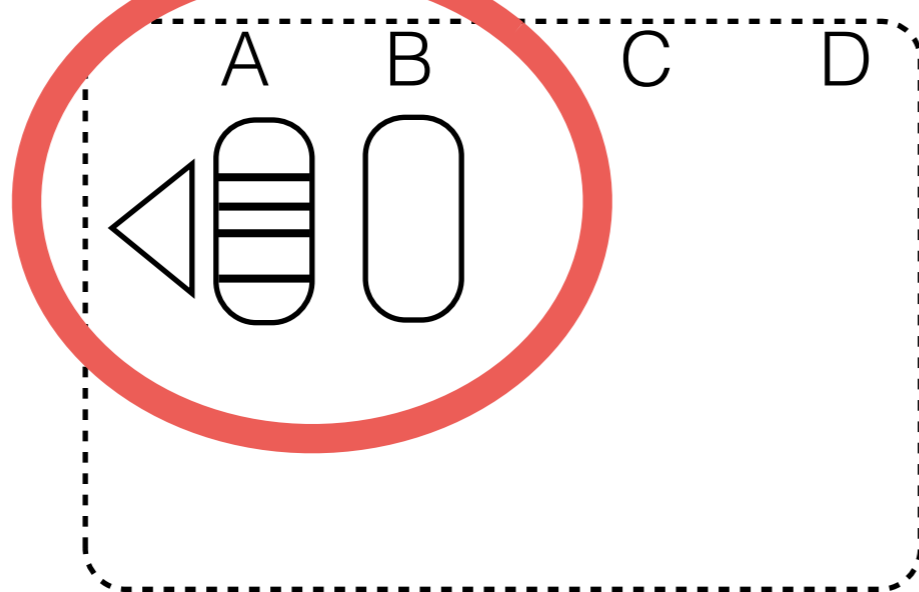


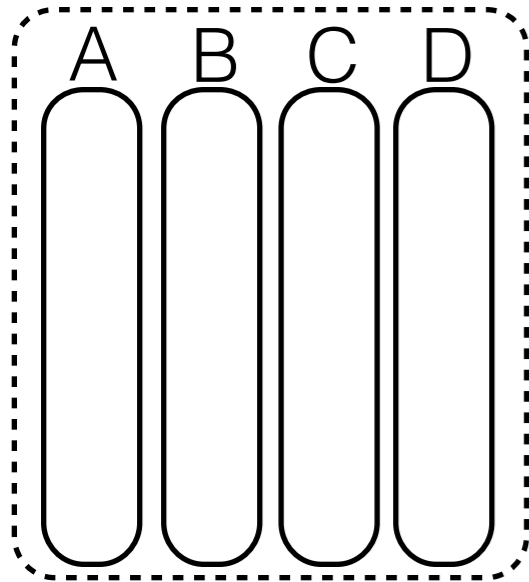
table 2



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

cracking tangram

base data
table 1



as queries arrive...
table 1

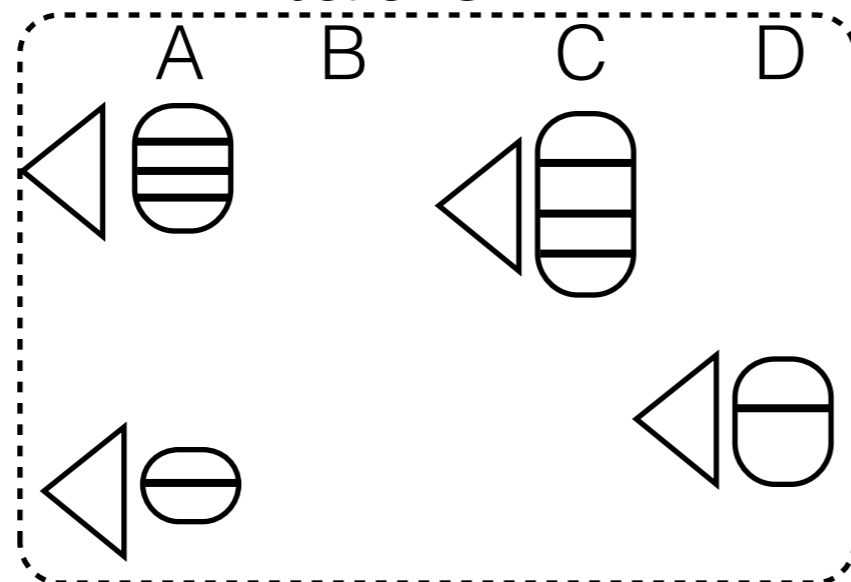


table 2

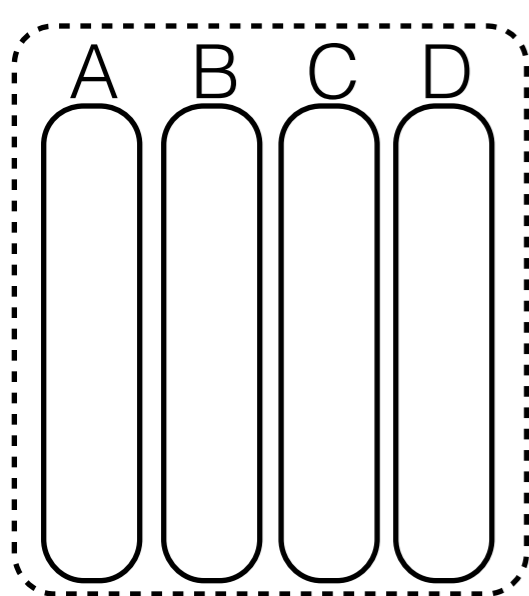
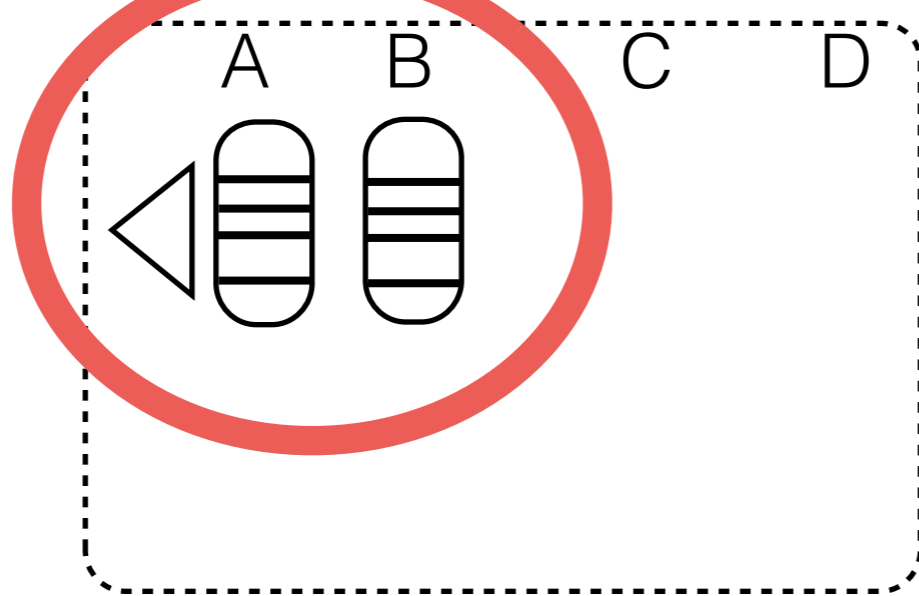


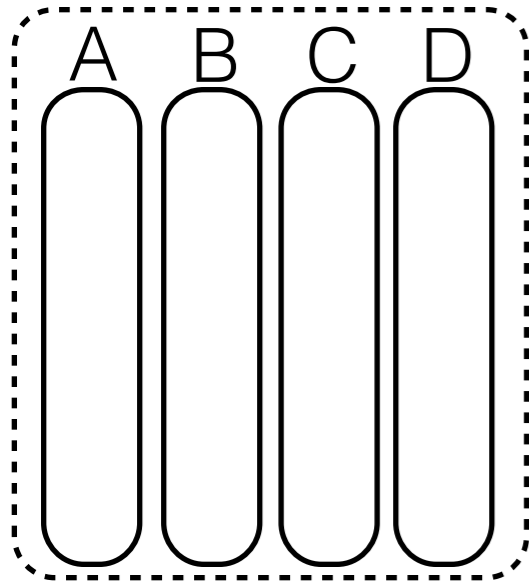
table 2



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

cracking tangram

base data
table 1



as queries arrive...
table 1

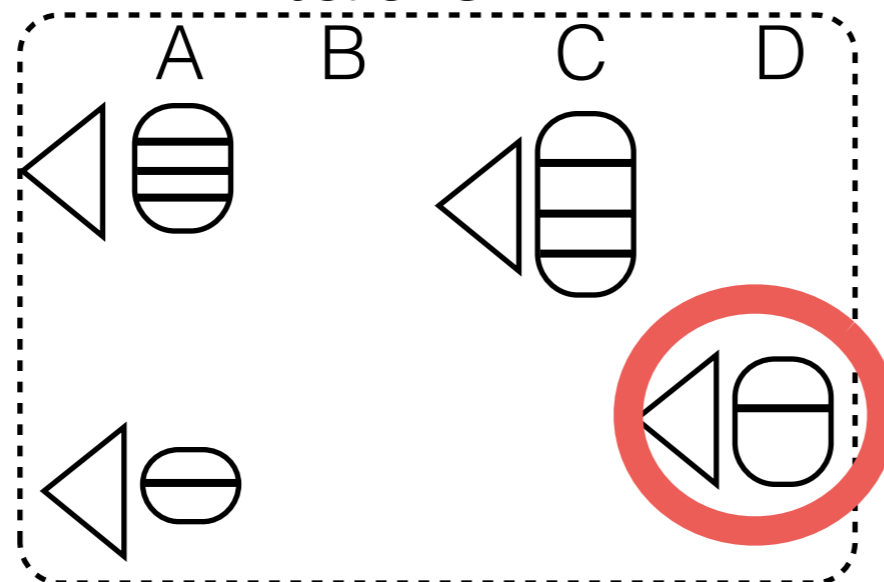


table 2

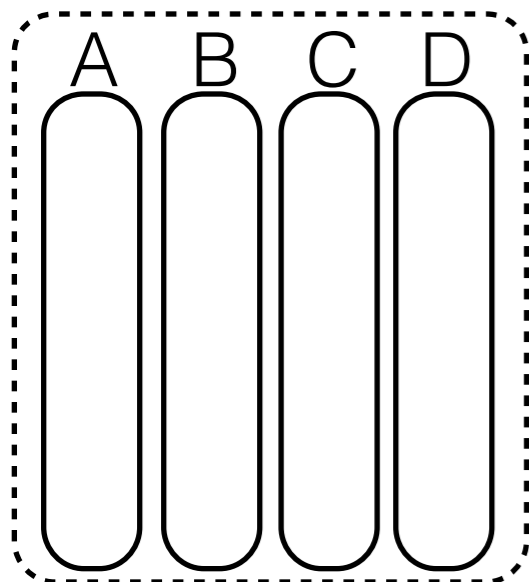
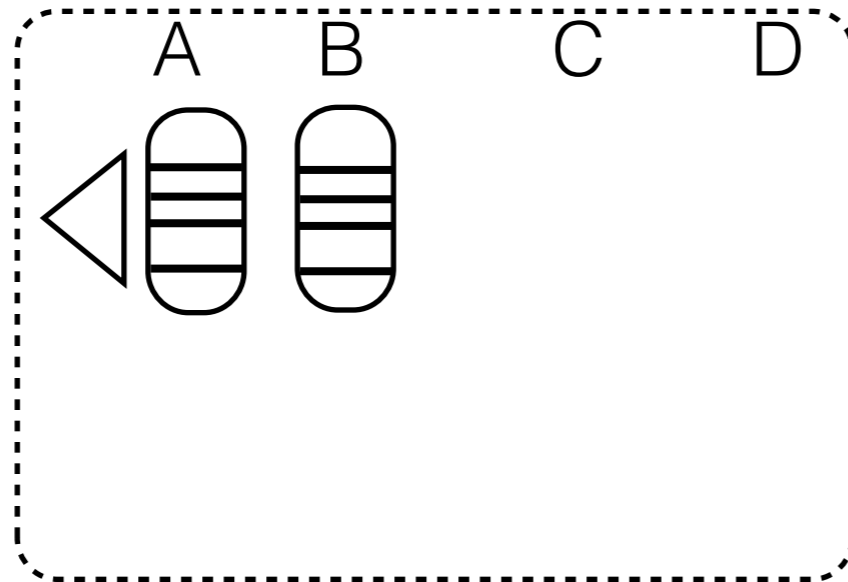


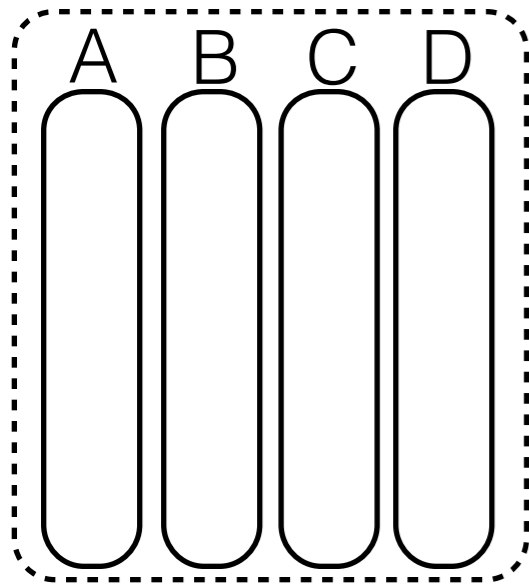
table 2



- 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

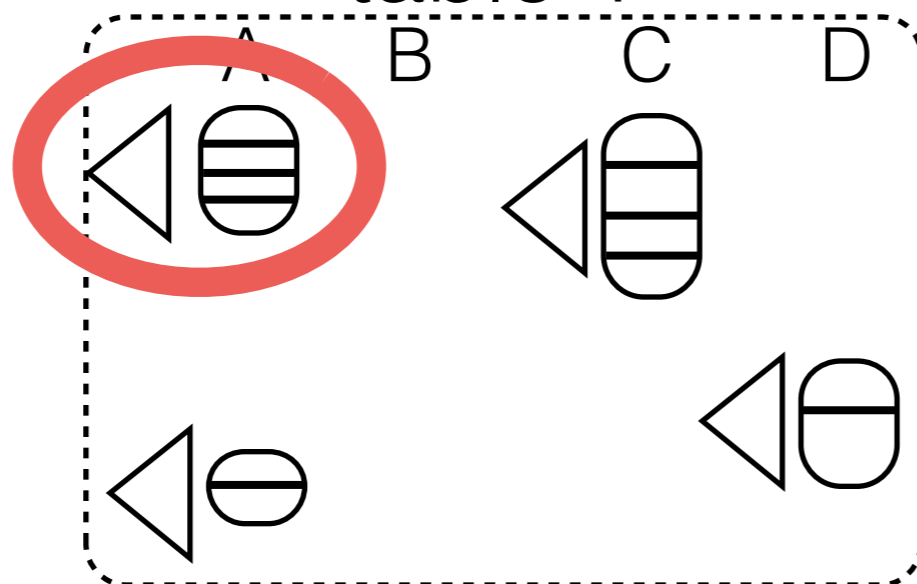


table 2

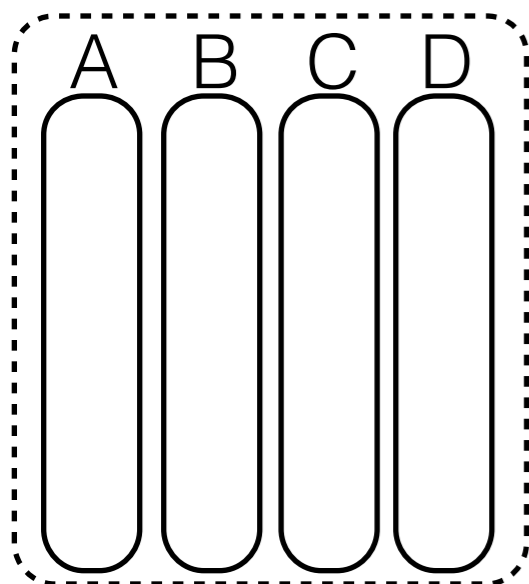
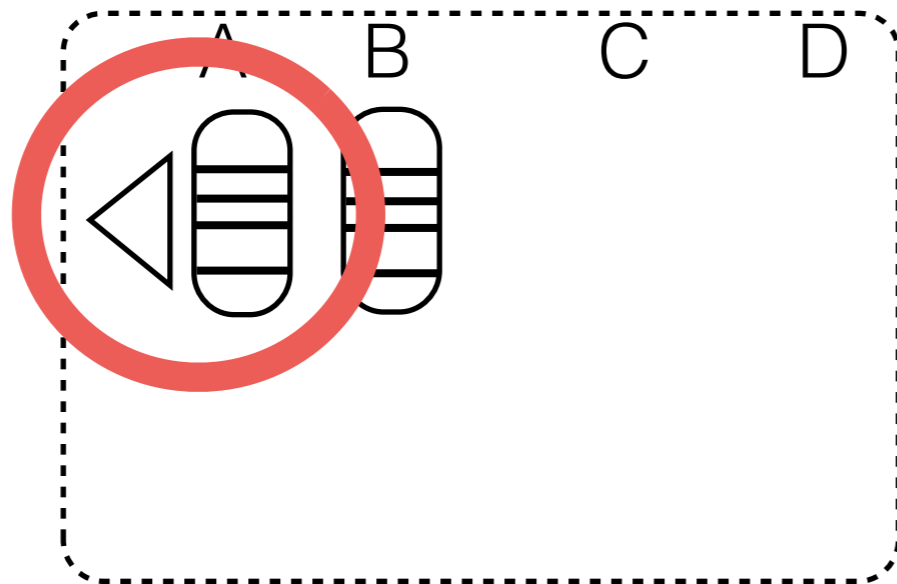


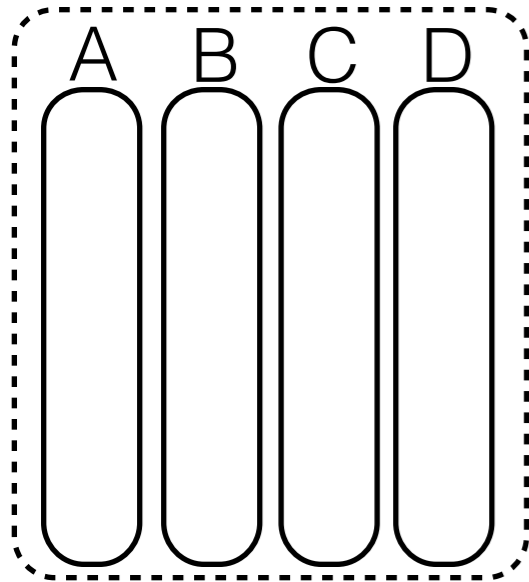
table 2



- 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

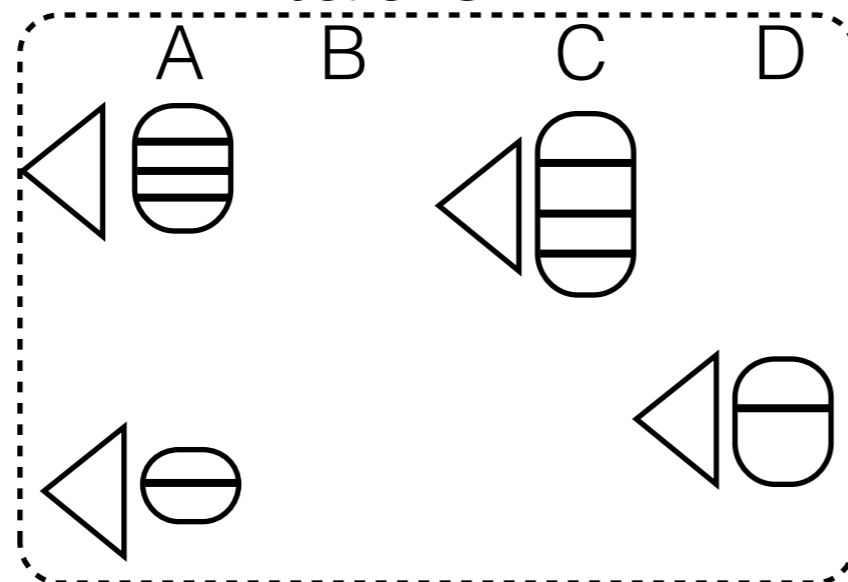


table 2

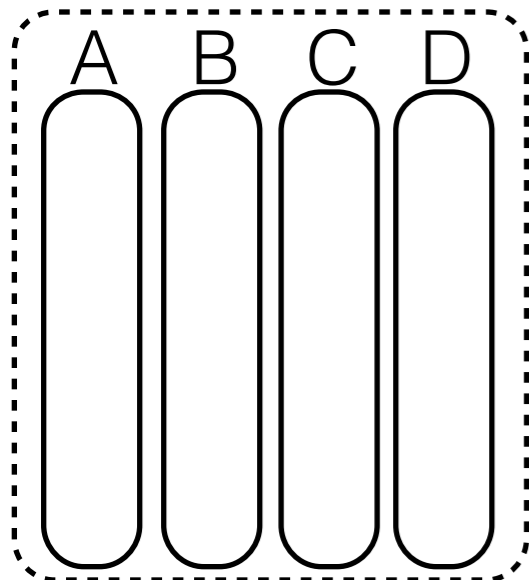
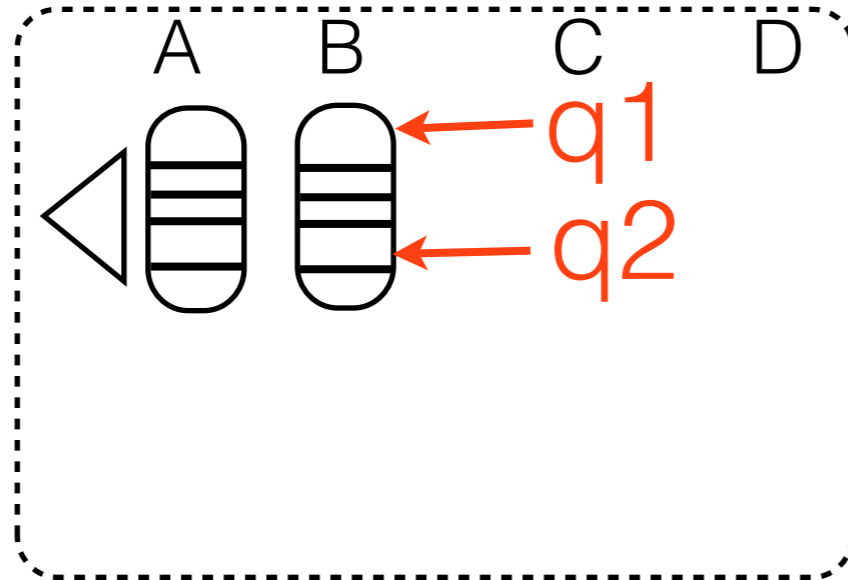


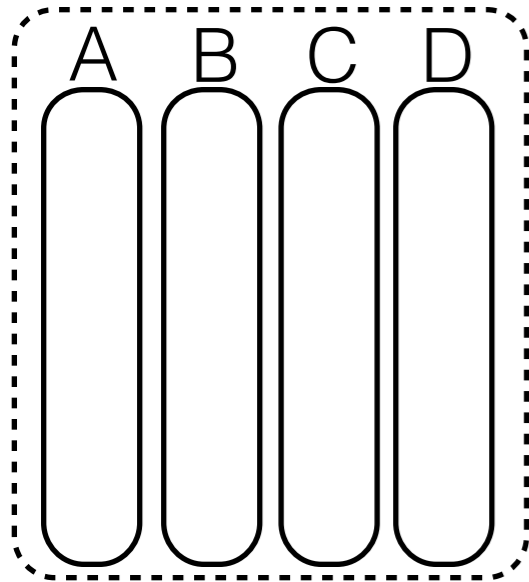
table 2



- 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

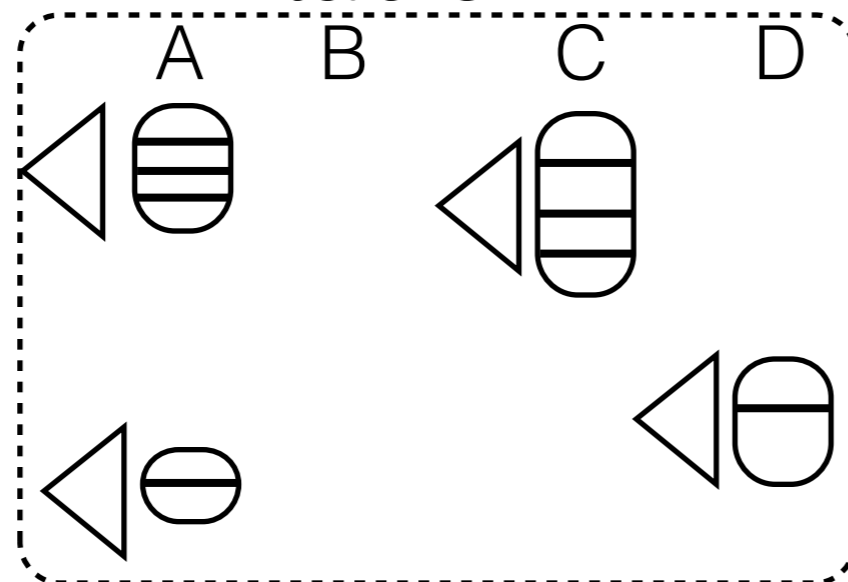


table 2

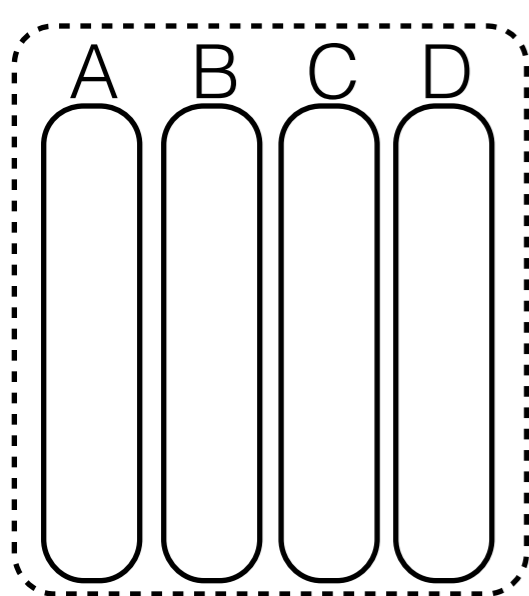
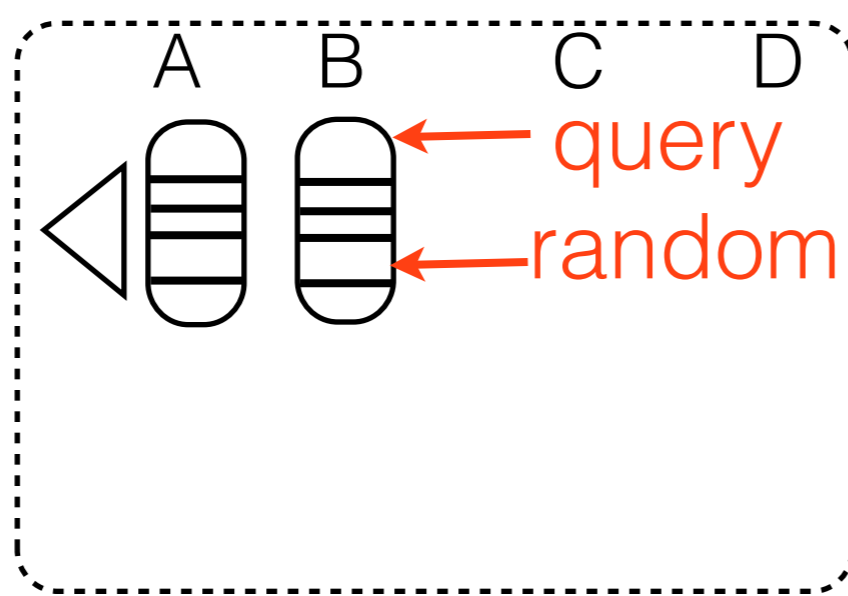
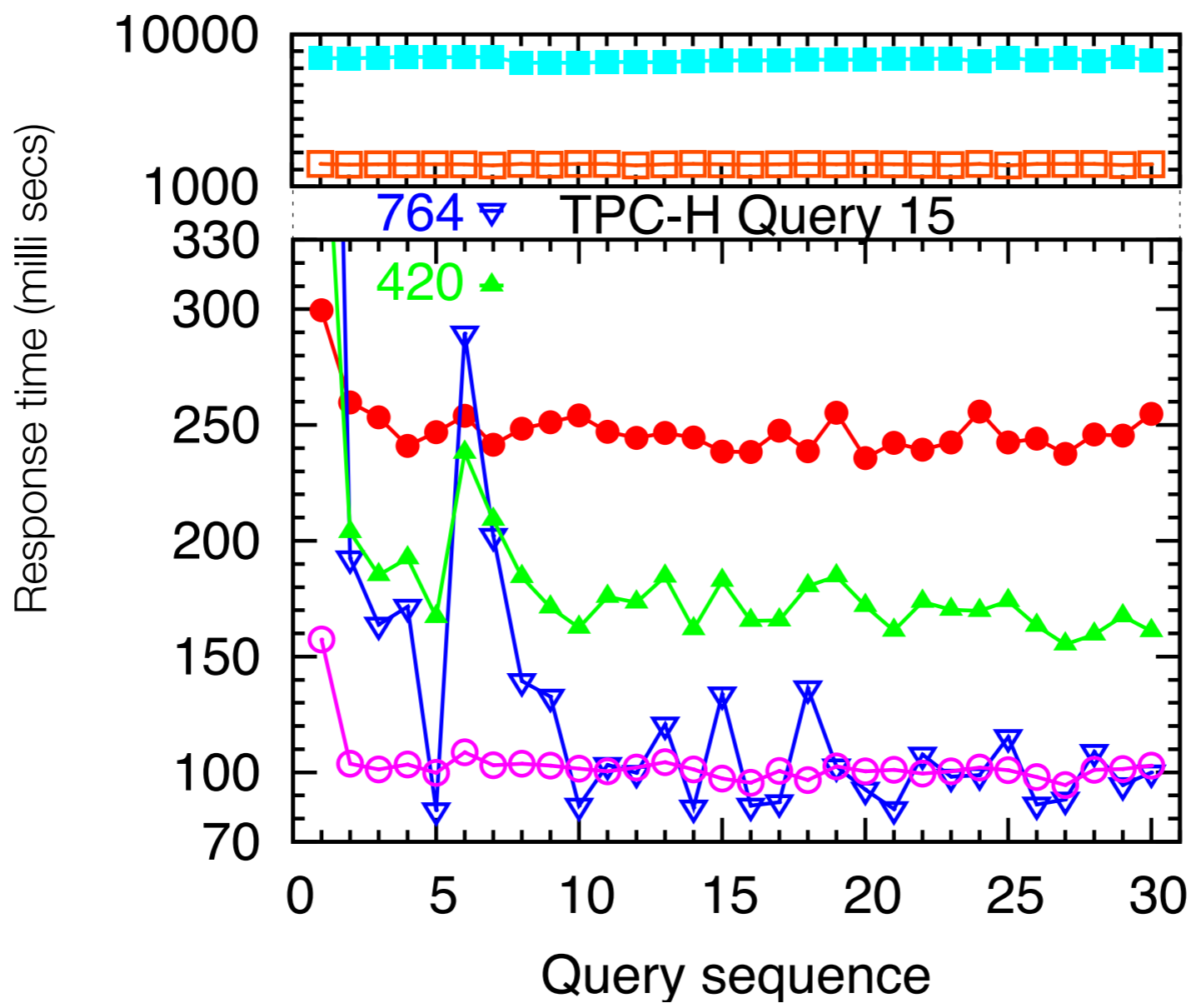
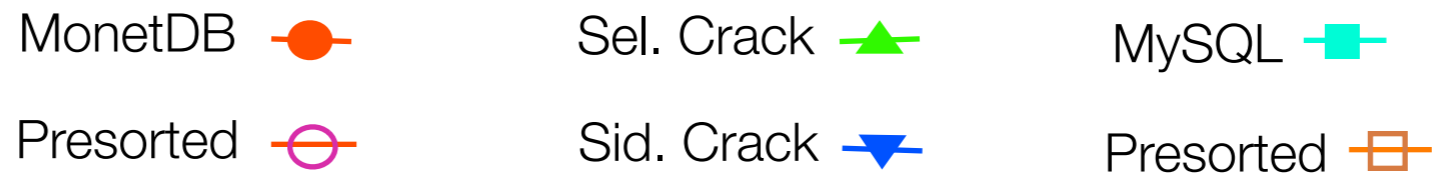
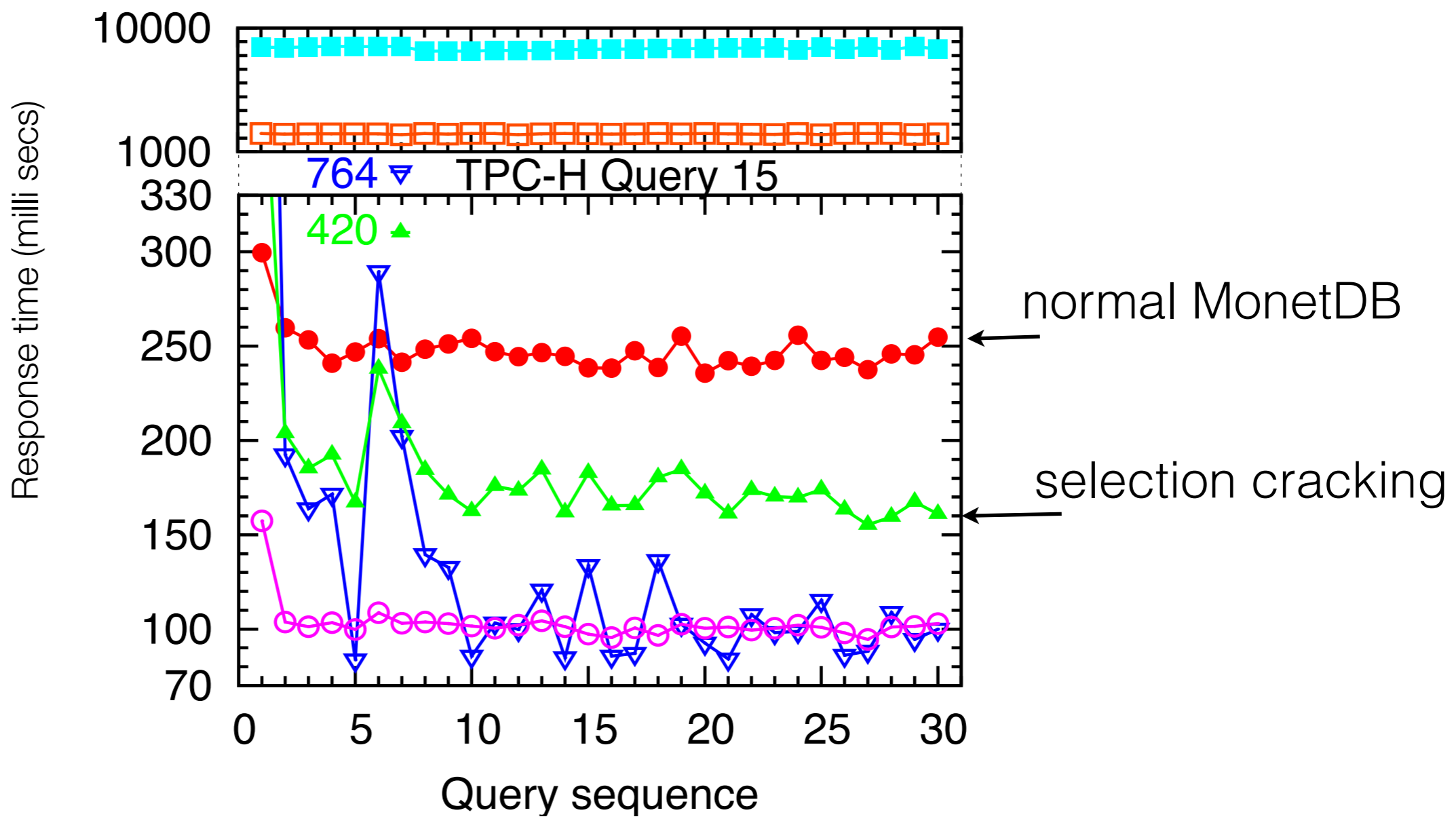
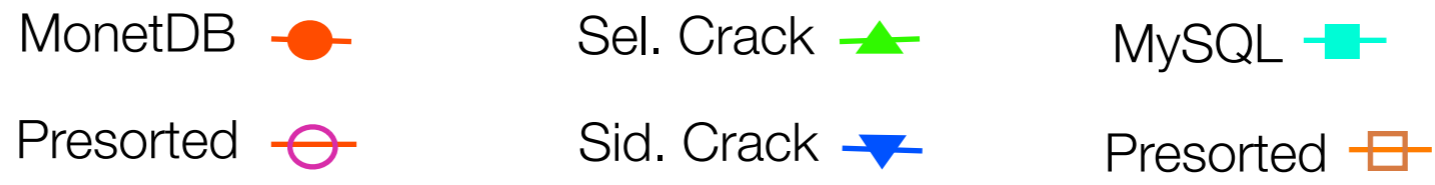


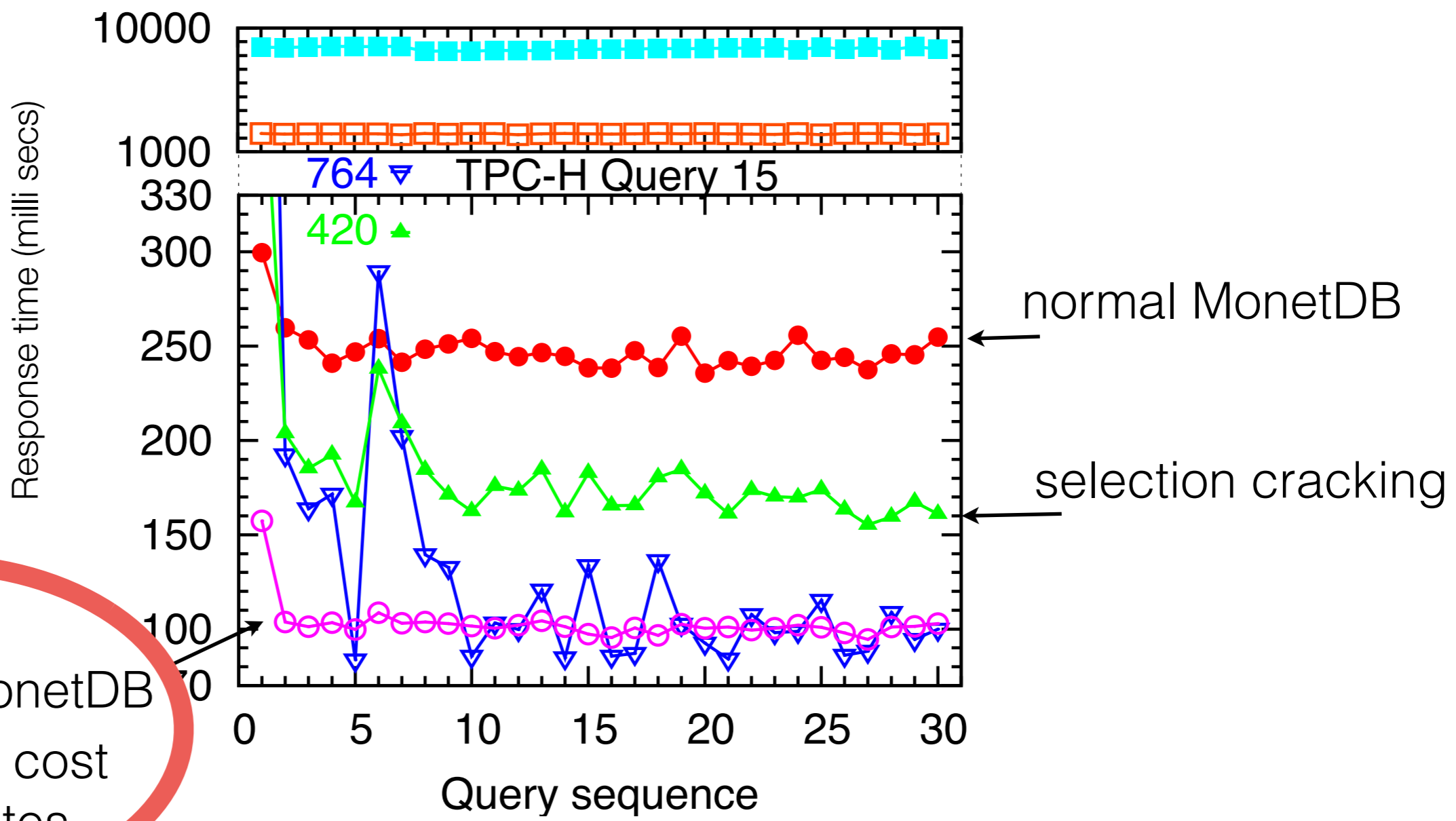
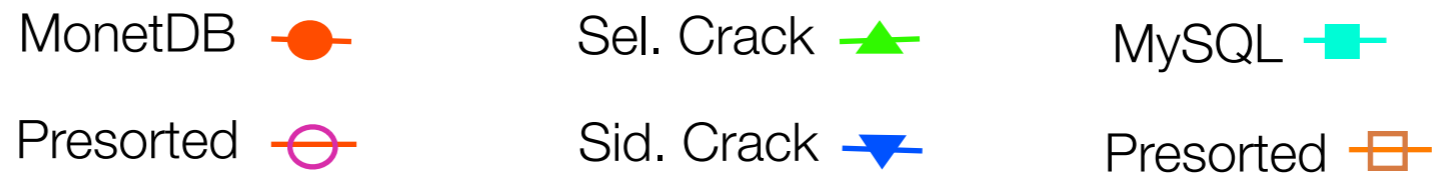
table 2



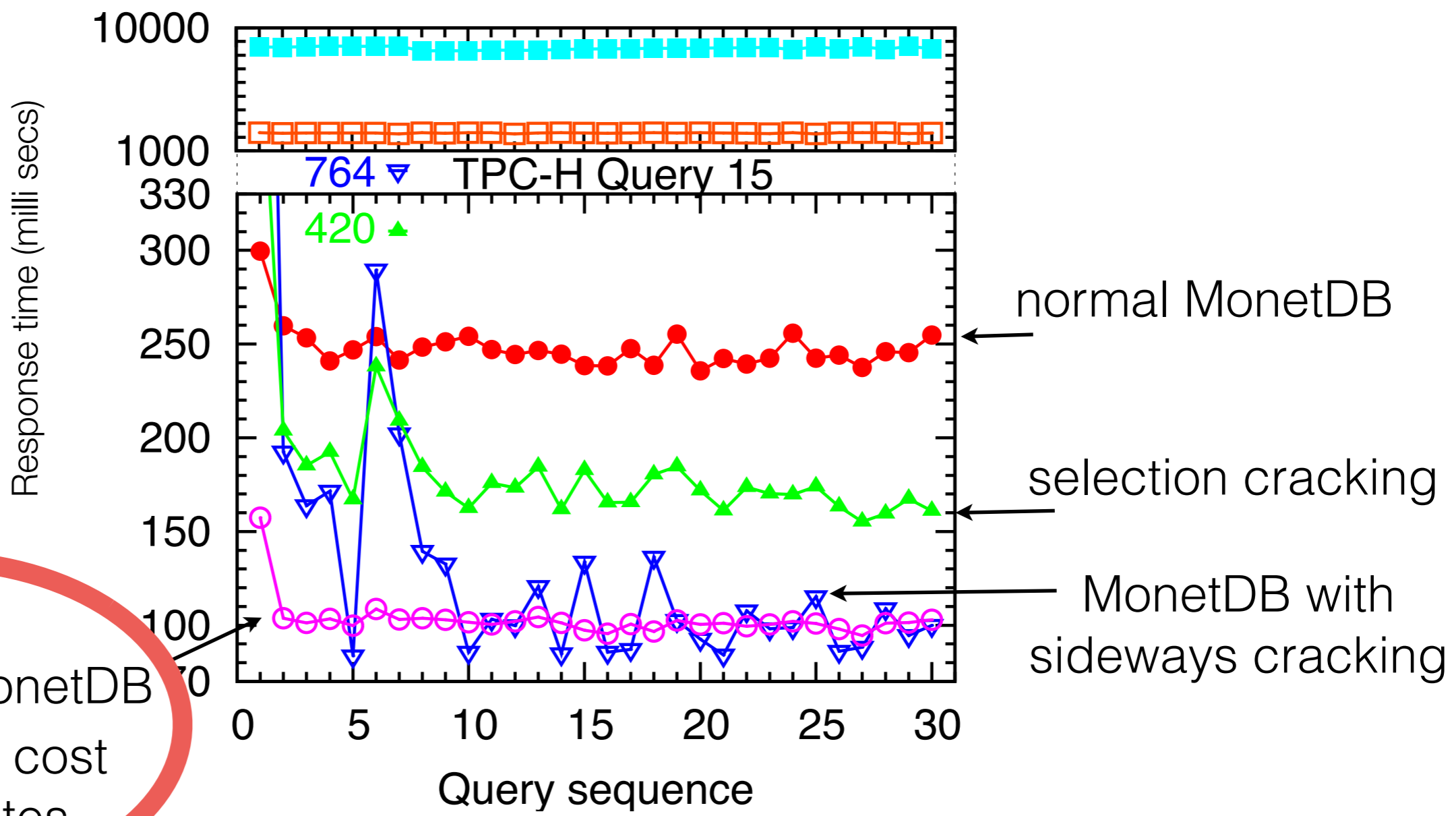
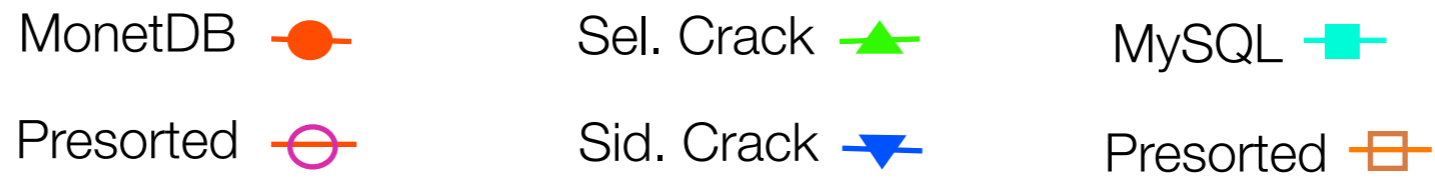
- partial materialization
- partial indexing
- continuous adaptation
- storage adaptation
- no tuple reconstruction
- adaptive alignment
- sort in caches
- crack joins
- lightweight locking
- stochastic cracking



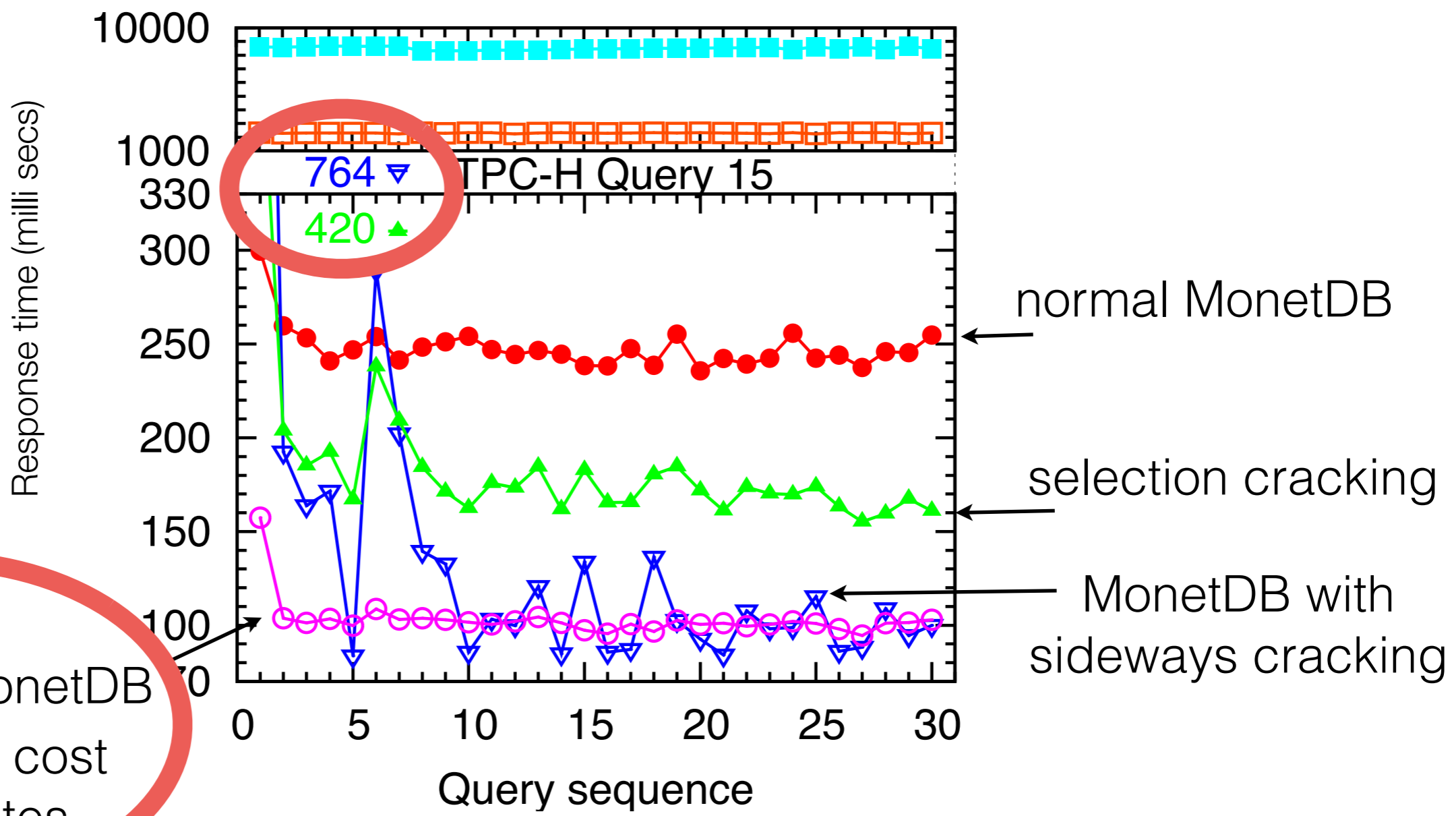
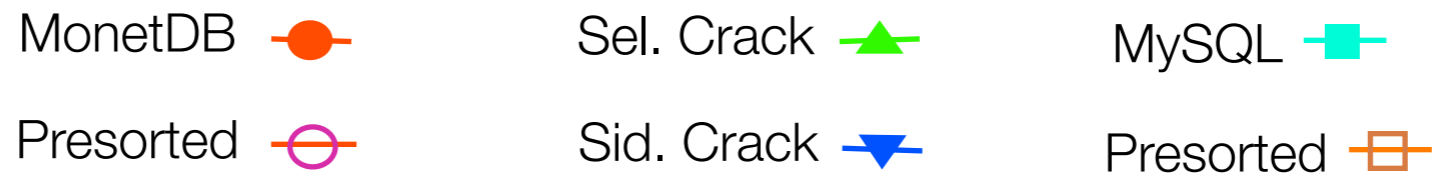


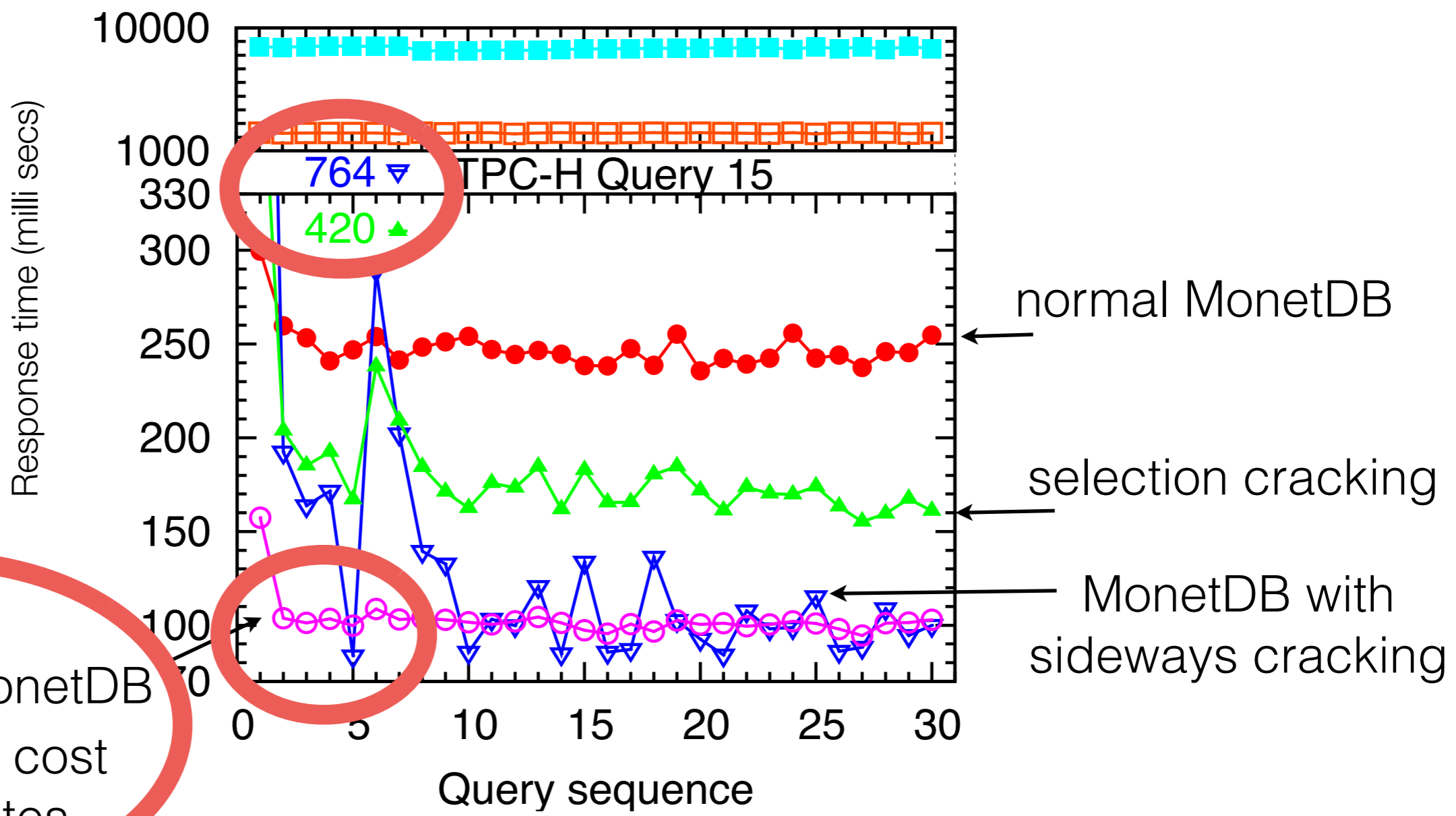
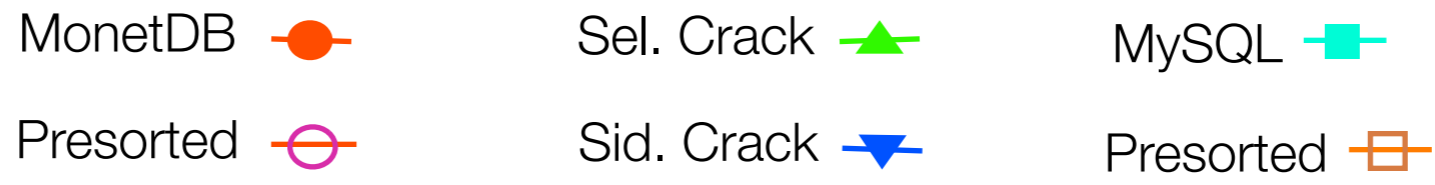


presorted MonetDB
preparation cost
3-14 minutes

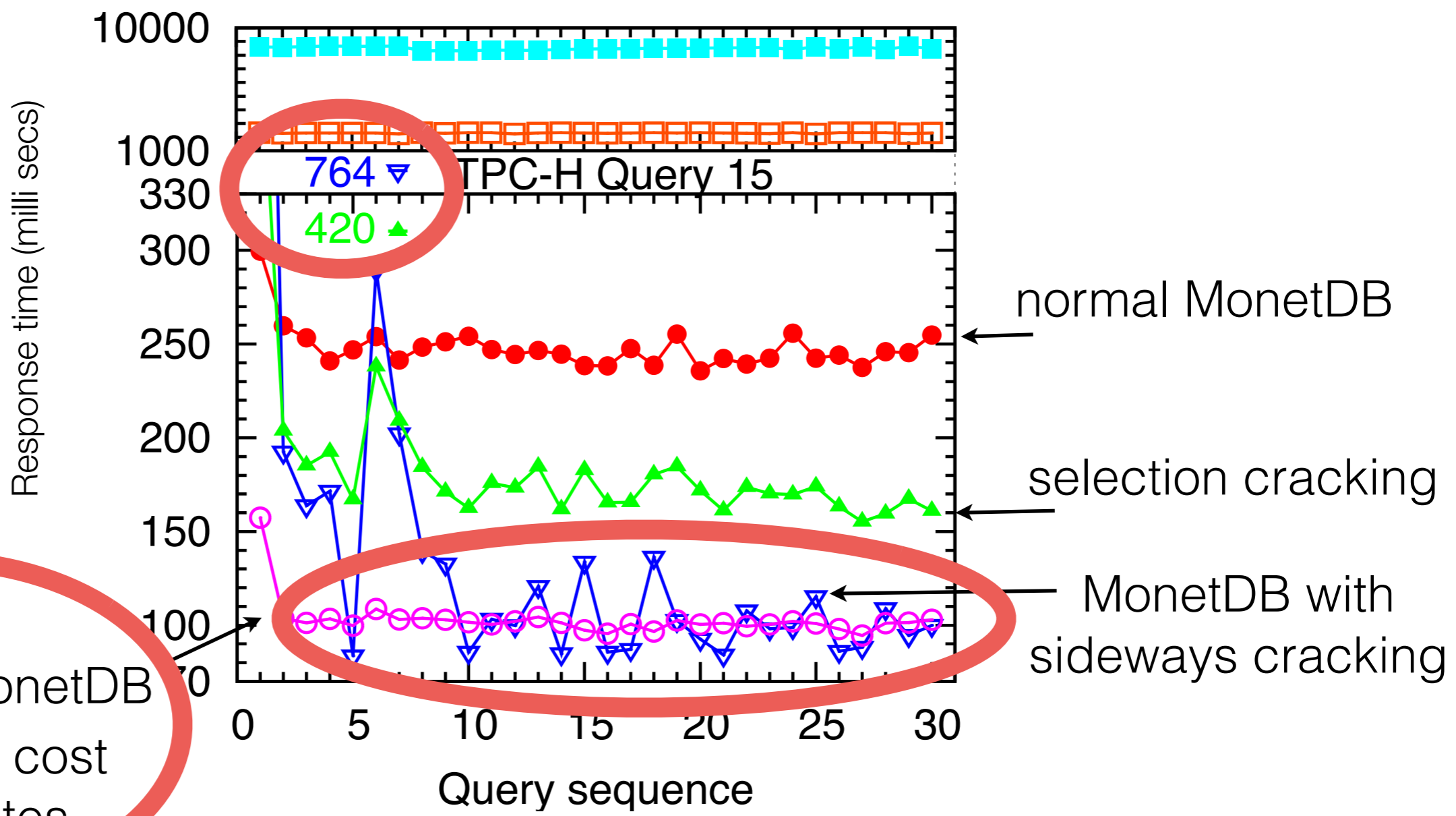
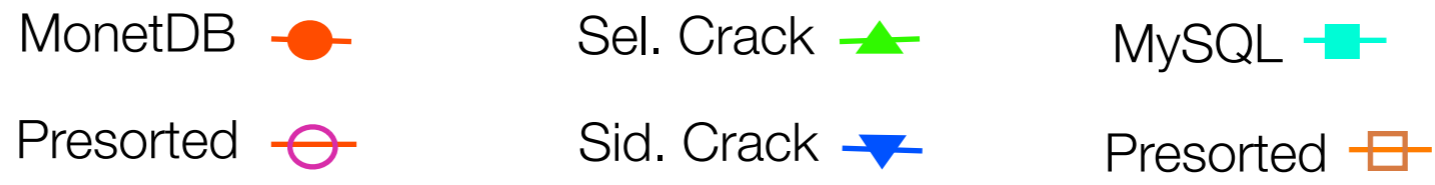


presorted MonetDB
preparation cost
3-14 minutes



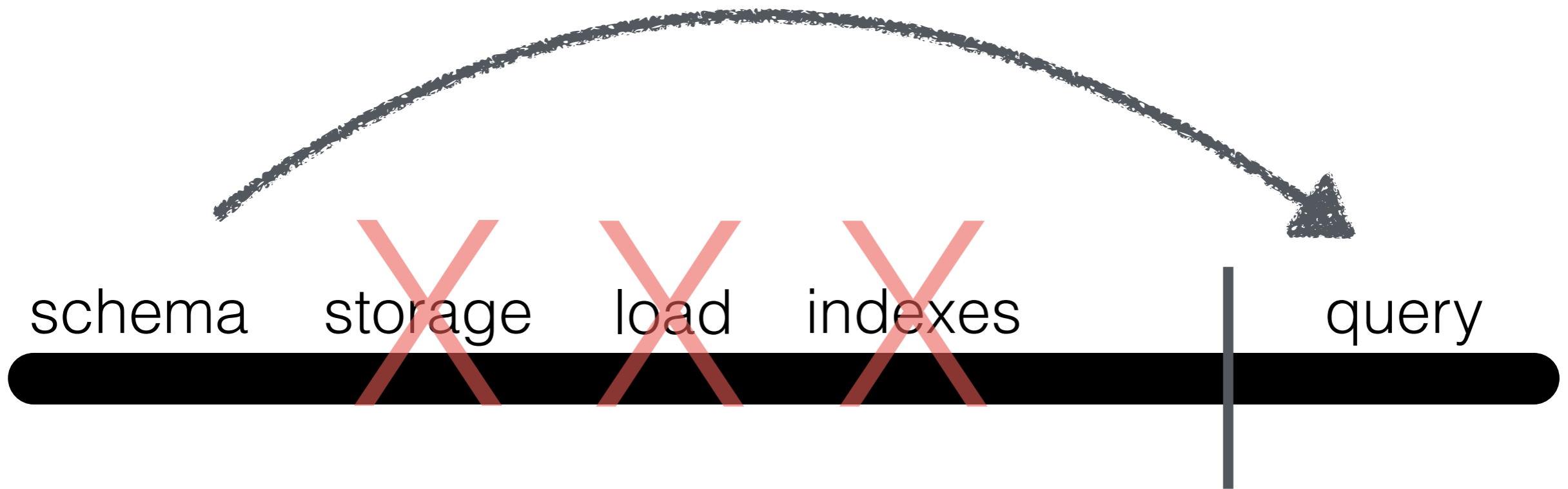


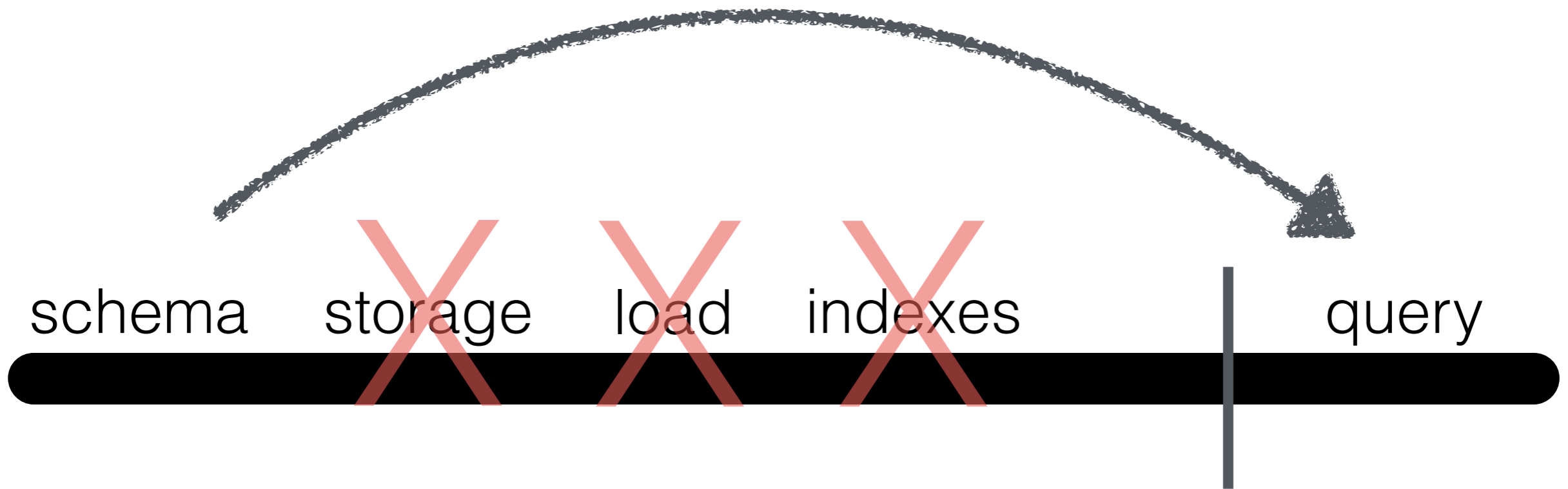
presorted MonetDB
preparation cost
3-14 minutes





cracking on Skyserver (4TB)
(Sloan Digital Sky Survey, www.sdss.org)

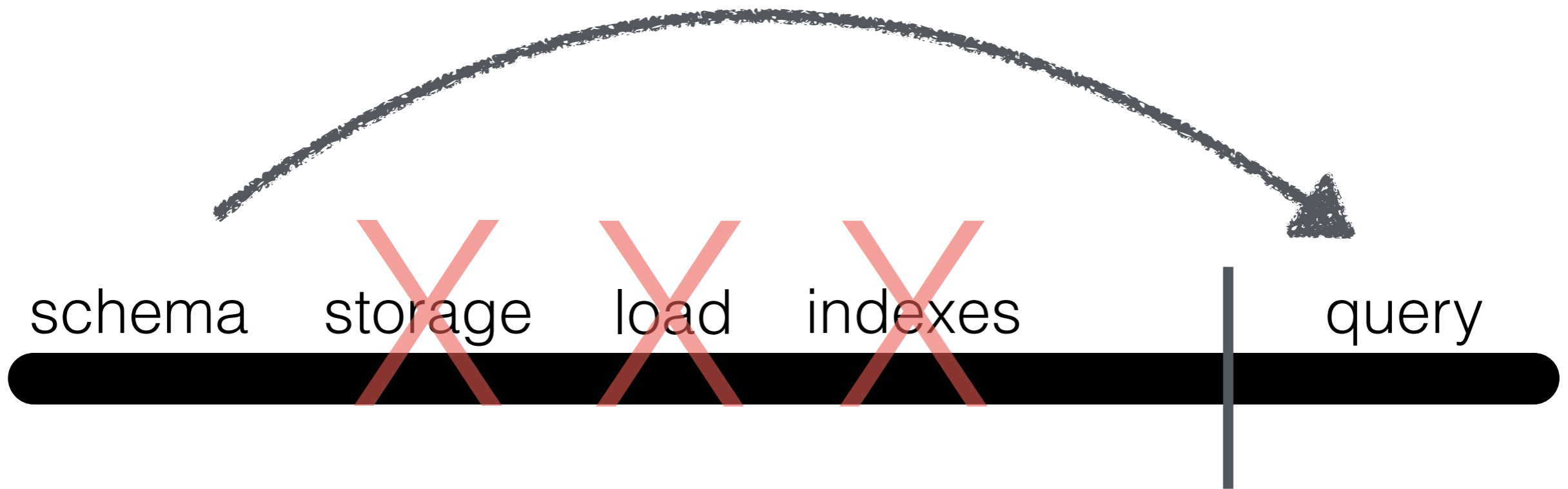
cracking answers 160.000 queries
while full indexing is still half way creating one index






adaptive loading (NoDB, CIDR11/SIGMOD12)

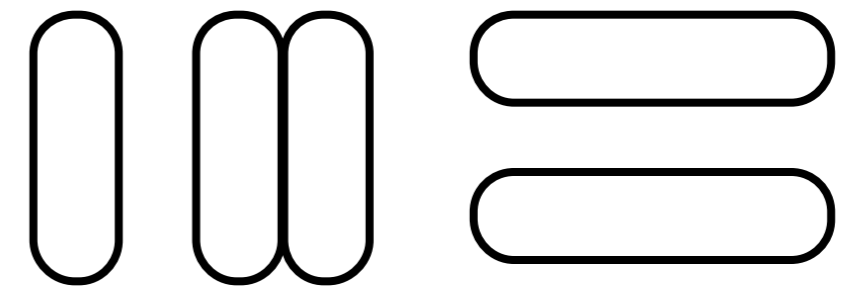
STILL 
LOADING 

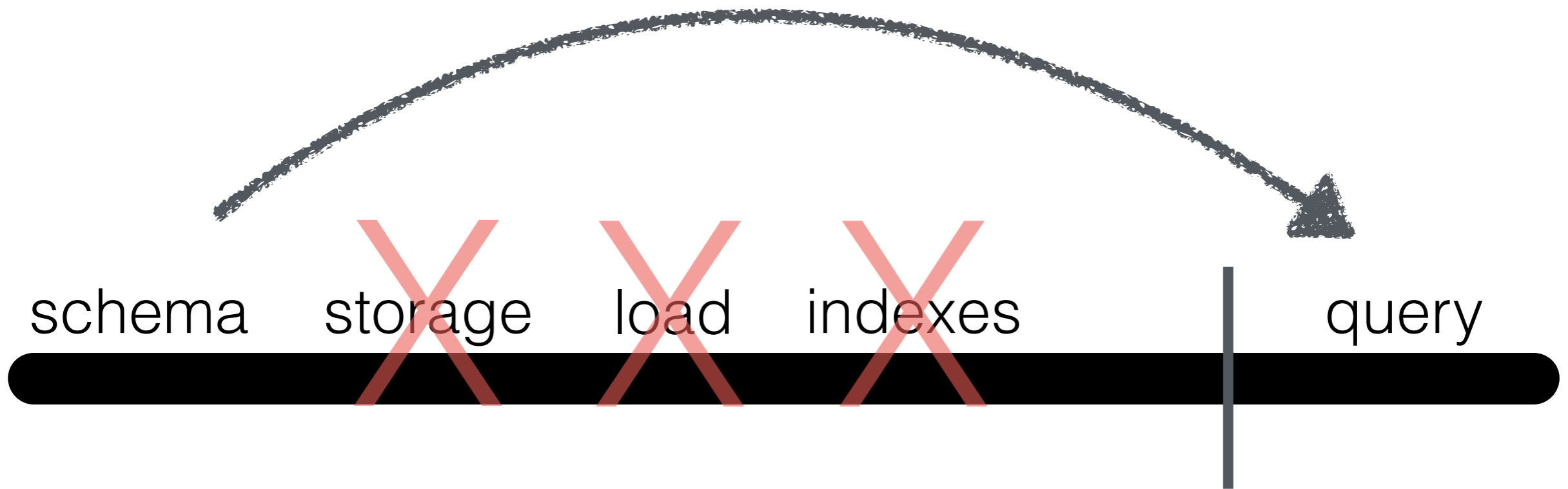


adaptive loading (NoDB, CIDR11/SIGMOD12)

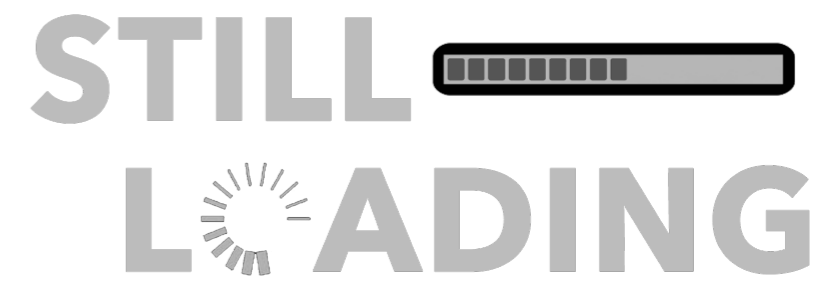
STILL 
LOADING

adaptive storage (H20, SIGMOD14)

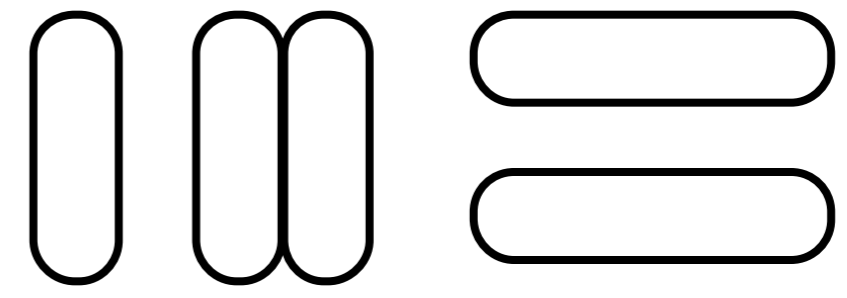




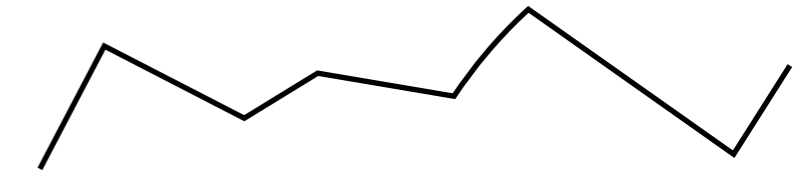
adaptive loading (NoDB, CIDR11/SIGMOD12)

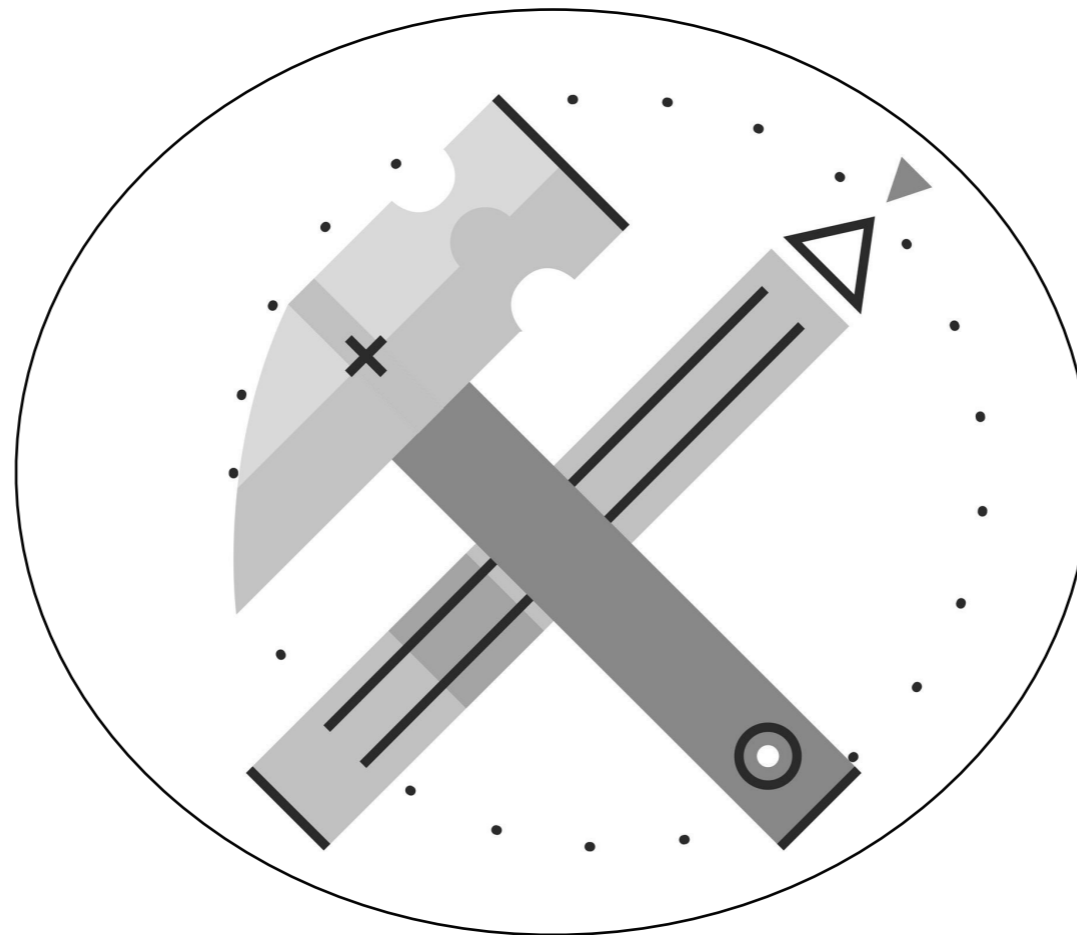


adaptive storage (H20, SIGMOD14)

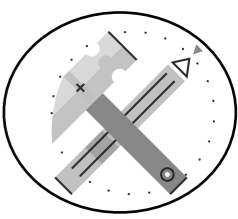


adaptive time series indexing (ADS, SIGMOD14)





data systems that are easy to design
(storage, data flow, algorithms, tuning, etc)



+



+



=



One size does
not fit all

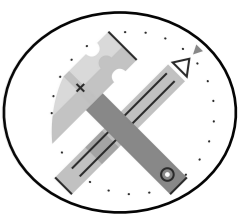
Custom solutions are needed
for optimal performance

Solutions need
to be tuned

Bootstrapping new systems is
expensive and time-consuming

e.g., column-stores:

first ideas in 80s,
first advanced architectures in 90s,
first rather complete designs in early 2000s,
industry adoption 2010+
still no indexing, cost based optimizations, ...



conflicting goals

moving target

(hardware and requirements change continuously and rapidly)



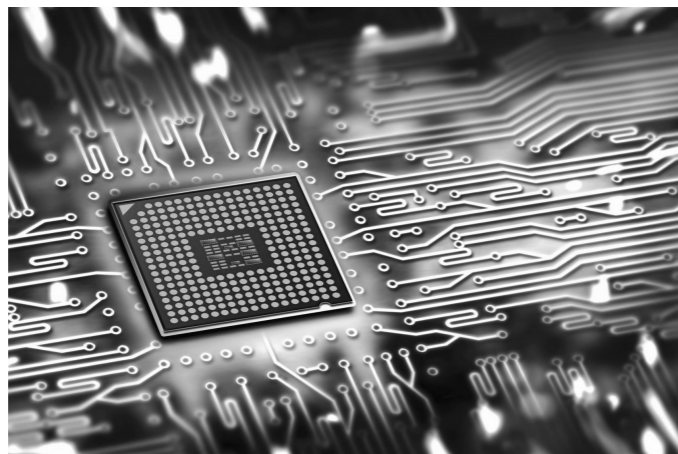
application requirements



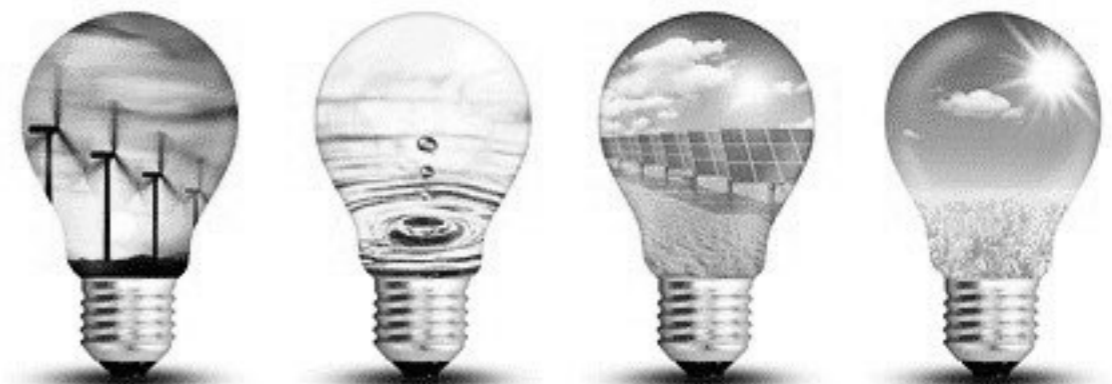
performance



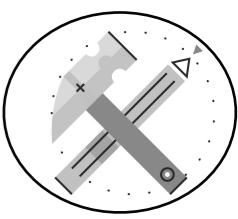
budget



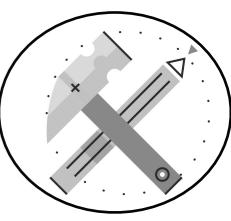
hardware



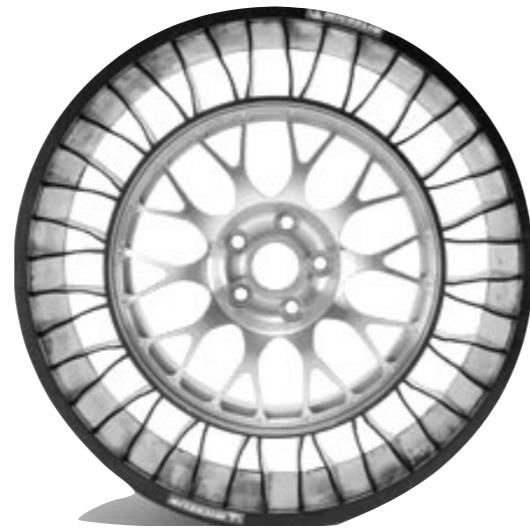
energy profile



data systems design (and research) is kind of an **art**



disk

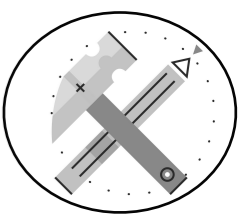


memory



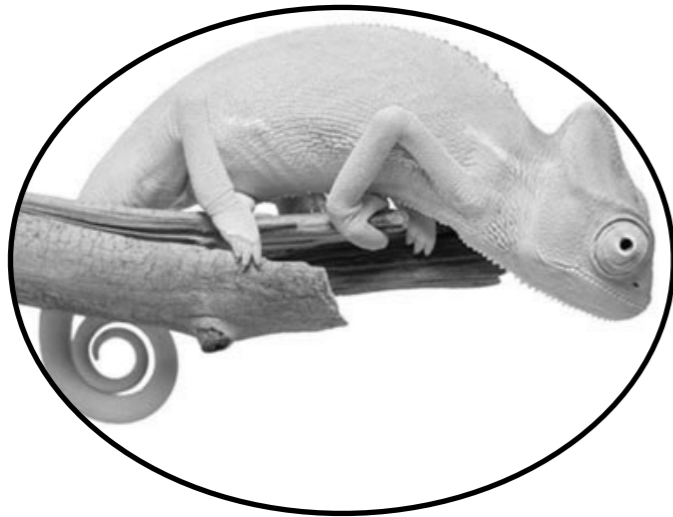
flash

...

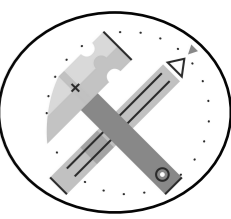


self-designing data systems

data+queries+hardware

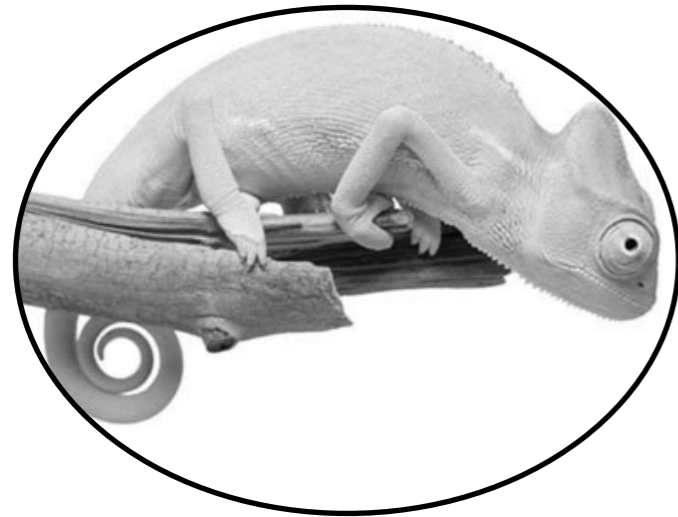


data system



self-designing data systems

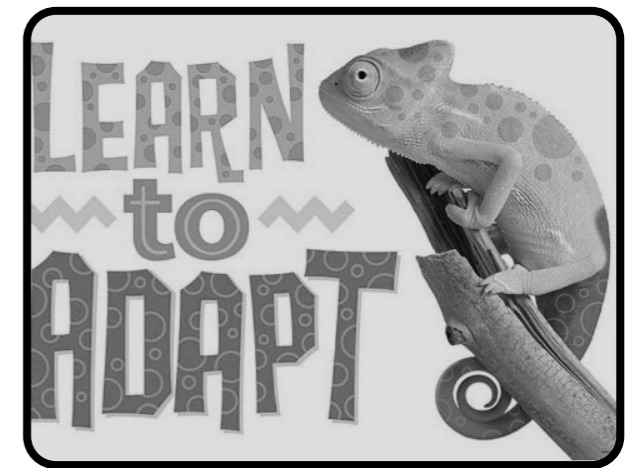
data+queries+hardware



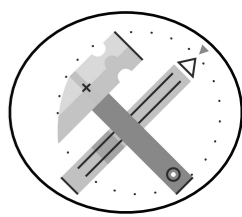
data system



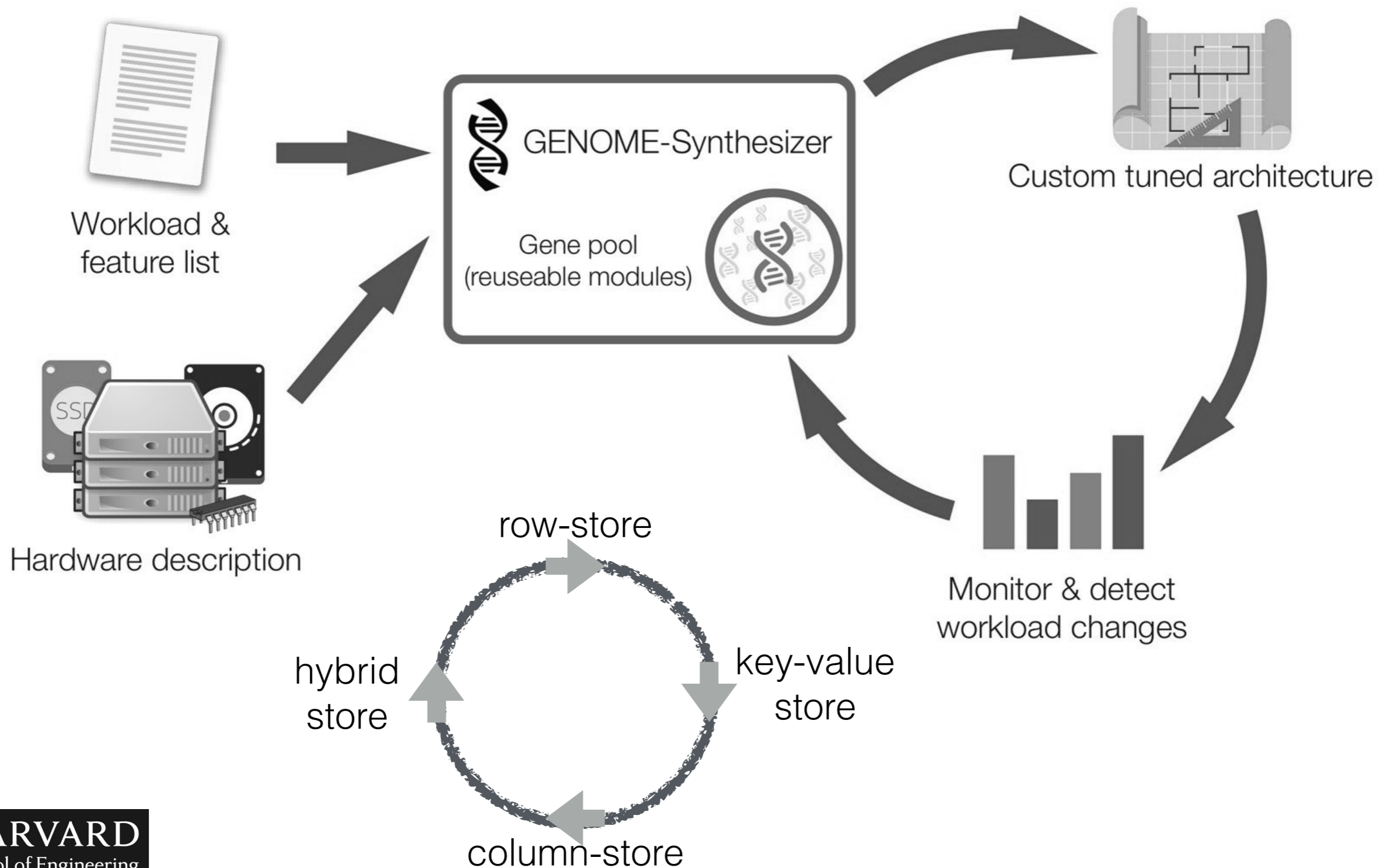
easy to design



adapt to environment



adaptivity across architecture borders





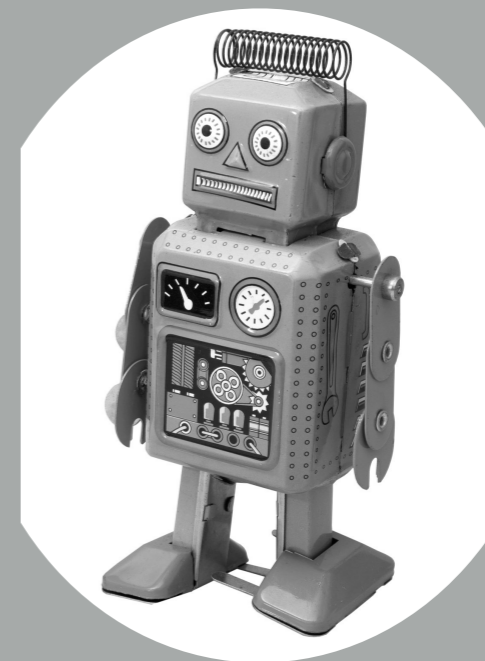
data systems that
are easy to use



dbTouch



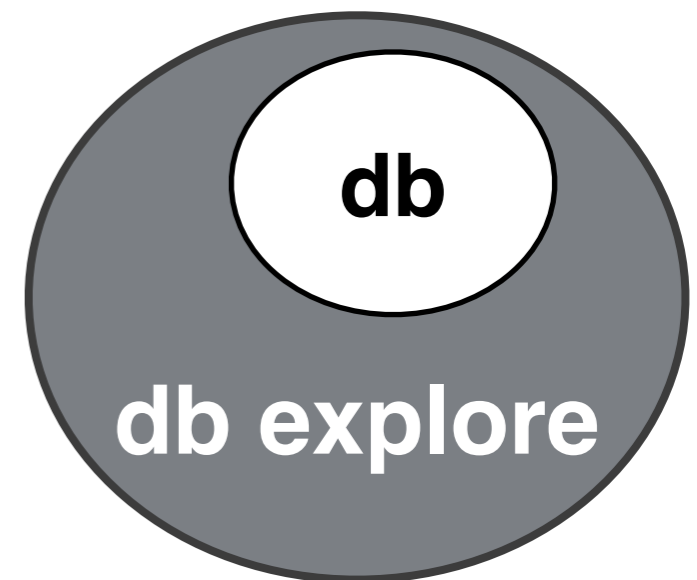
show me something
interesting



— DATA

Queriosity

data systems today
allow us to answer queries fast



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

instead of making fixed decisions

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

DATA SYSTEMS LABORATORY

@ Harvard School of Engineering and Applied Sciences

Designing data systems for the big data era

<http://daslab.seas.harvard.edu/>

Martin Kersten
Stefan Manegold
Goetz Graefe
Harumi Kuno
Anastasia Ailamaki
Themis Palpanas
Eleni Petraki

Ioannis Alagiannis
Miguel Branco
Renata Borovica
Erietta Liarou
Felix Halim
Ronald Yap
Panos Karras

Kostas Zoumpatianos
Manos Athanassoulis
Lukas Maas
Abdul Wasay
Mike Kester
Dhruv Gupta

DATA SYSTEMS LABORATORY

@ Harvard School of Engineering and Applied Sciences

Designing data systems for the big data era

<http://daslab.seas.harvard.edu/>

thank you!

Martin Kersten
Stefan Manegold
Goetz Graefe
Harumi Kuno
Anastasia Ailamaki
Themis Palpanas
Eleni Petraki

Ioannis Alagiannis
Miguel Branco
Renata Borovica
Erietta Liarou
Felix Halim
Ronald Yap
Panos Karras

Kostas Zoumpatianos
Manos Athanassoulis
Lukas Maas
Abdul Wasay
Mike Kester
Dhruv Gupta