

Modern Dataflow

Frank McSherry

@frankmcsherry

<http://github.com/frankmcsherry>

Acks: Derek Murray, Rebecca Isaacs, Michael Isard, Paul Barham, Martin Abadi, Gordon Plotkin, Andrea Lattuada, Zaheer Chothia, John Liagouris,

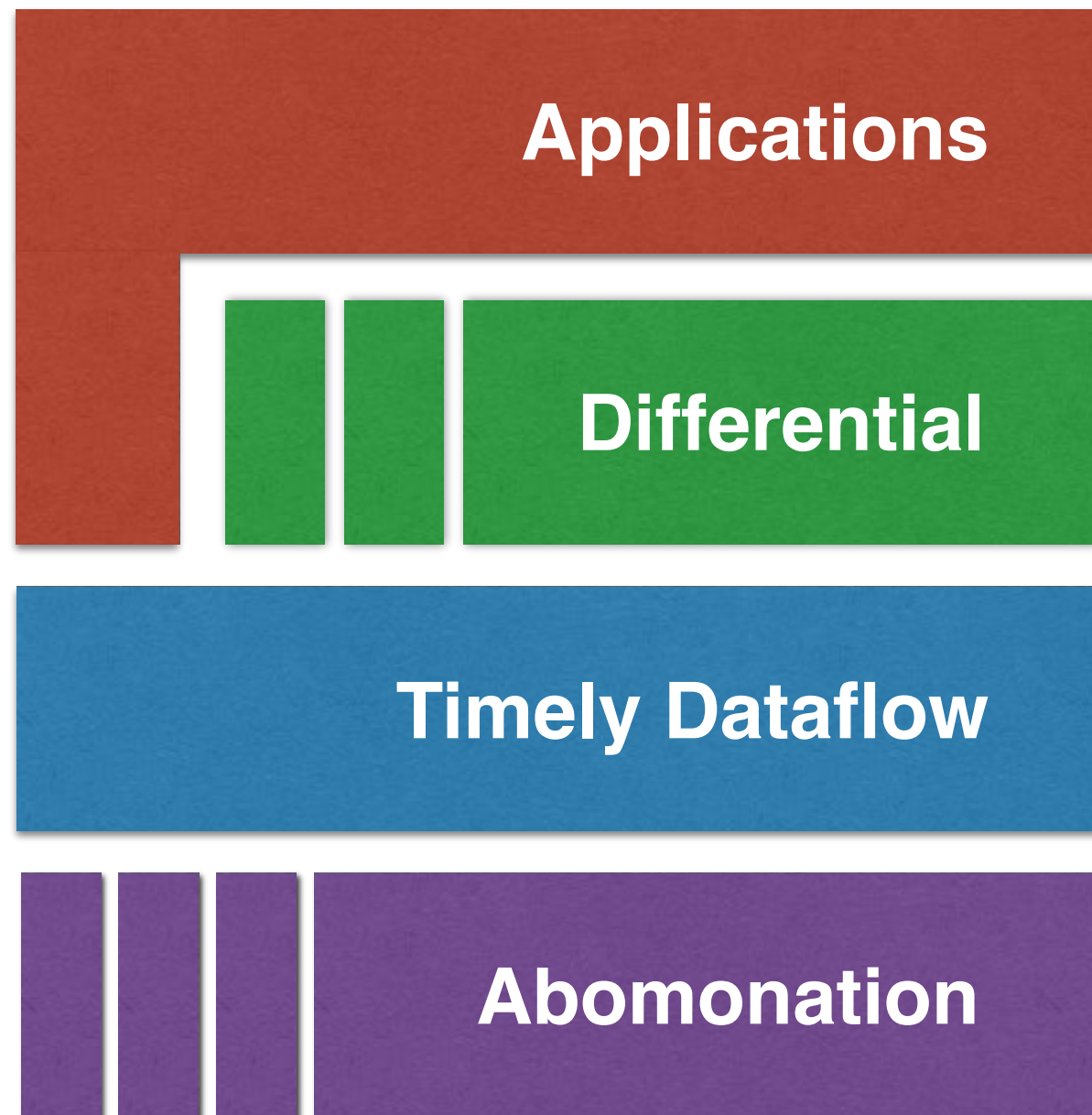
20xPR	cores	twitter_rv	uk_2007_05
Spark	128	857s	1759s
Giraph	128	596s	1235s
GraphLab	128	249s	833s
GraphX	128	419s	462s

[GraphX, OSDI 2014]

20xPR	cores	twitter_rv	uk_2007_05
Spark	128	857s	1759s
Giraph	128	596s	1235s
GraphLab	128	249s	833s
GraphX	128	419s	462s
Laptop	1	110s	256s

[COST, HotOS 2015]

20xPR	cores	twitter_rv	uk_2007_05
Spark	128	857s	1759s
Giraph	128	596s	1235s
GraphLab	128	249s	833s
GraphX	128	419s	462s
Laptop	1	110s	256s
Timely	128	15s	19s



Cool stuff like real-time graph analytics, etc.

Language for scalable incremental computation

A bit like an OS for data-parallel compute

Data serialization at memory bandwidth

Differential Dataflow

“Functional Reactive Programming at Scale”

<http://github.com/frankmcsherry/differential-dataflow>

[based off of Differential Dataflow, CIDR 2013]

People are good at programming with collections.
(at least, better than with streams)

Gist: collection-oriented programming language,
the system manages changes to collections.

```
fn your_prog: [D] -> [R] = /* .. */;  
  
for t in times {  
    let output[t] = your_prog(input[t]);  
}
```

d_output: Stream<(Data, Time, isize)>

input streams of changes

```
let nodes = /* pairs (node, bool) */;  
let edges = /* pairs (node, node) */;
```

“program” : dataflow assembly

```
nodes.join(edges)           // one hop neighbors  
    .concat(nodes)         // plus original nodes  
    .distinct()            // extended neighborhood
```

dataflow execution

```
for t in times {  
    nodes.insert(..);  
    edges.insert(..);  
}
```



```
let nodes = /* pairs (node, bool) */;  
let edges = /* pairs (node, node) */;
```

```
nodes.join(edges)      // one hop neighbors  
    .concat(nodes)     // plus original nodes  
    .distinct()        // extended neighborhood
```

```
for t in times {  
    nodes.insert(..);  
    edges.insert(..);  
}
```

```
let nodes = /* pairs (node, bool) */;  
let edges = /* pairs (node, node) */;
```

```
nodes.join(edges)      // one hop neighbors  
    .concat(nodes)     // plus original nodes  
    .distinct()        // extended neighborhood
```

```
for t in times {  
    nodes.insert(..); nodes.remove(..);  
    edges.insert(..); edges.remove(..);  
}
```

```
let nodes = /* pairs (node, bool) */;  
let edges = /* pairs (node, node) */;  
  
    nodes.join(edges)           // one hop neighbors  
        .concat(nodes)         // plus original nodes  
        .distinct()             // extended neighborhood  
  
for t in times {  
    nodes.insert(..); nodes.remove(..);  
    edges.insert(..); edges.remove(..);  
}
```

```
let nodes = /* pairs (node, bool) */;  
let edges = /* pairs (node, node) */;  
  
nodes.iterate(|reach| {  
    nodes.join(edges)           // one hop neighbors  
        .concat(nodes)         // plus original nodes  
        .distinct()            // extended neighborhood  
});  
  
for t in times {  
    nodes.insert(..); nodes.remove(..);  
    edges.insert(..); edges.remove(..);  
}
```

```
let nodes = /* pairs (node, bool) */;  
let edges = /* pairs (node, node) */;  
  
nodes.iterate(|reach| {  
    reach.join(edges) // one hop neighbors  
        .concat(nodes) // plus original nodes  
        .distinct() // extended neighborhood  
});  
    Stream<((node, bool), (Time, u64), isize)>  
  
for t in times {  
    nodes.insert(..); nodes.remove(..);  
    edges.insert(..); edges.remove(..);  
}
```

Reach	cores	livejournal	orkut	● ●
GraphX	128	36s	48s	
Socialite	128	52s	67s	● ●
Myria	128	5s	6s	
BigDatalog	128	17s	20s	● ●

[BigDatalog, SIGMOD 2016]

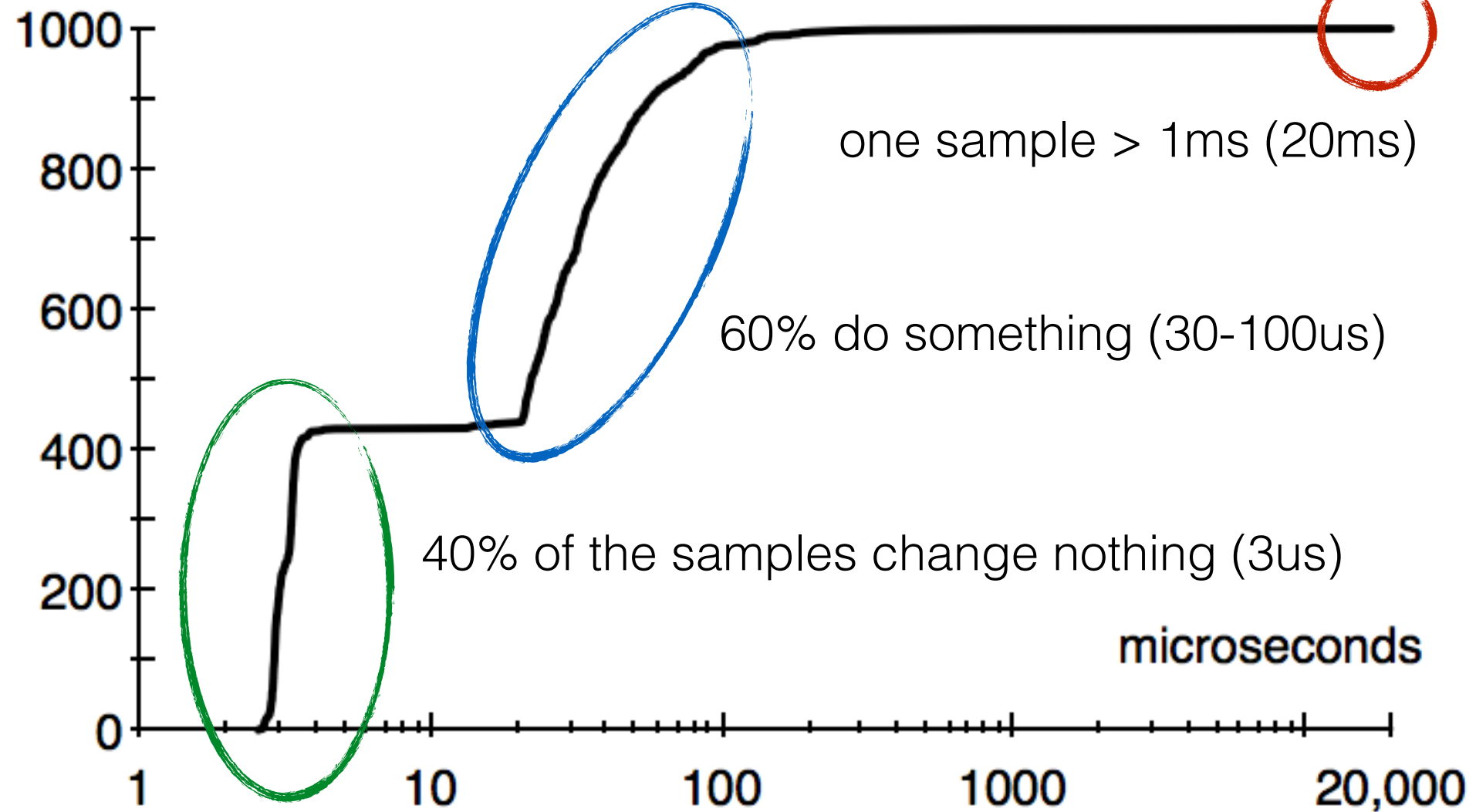
Reach	cores	livejournal	orkut
GraphX	128	36s	48s
Socialite	128	52s	67s
Myria	128	5s	6s
BigDatalog	128	17s	20s
Differential	1	7s	15s

Reach

cores

livejournal

orkut



update

50us

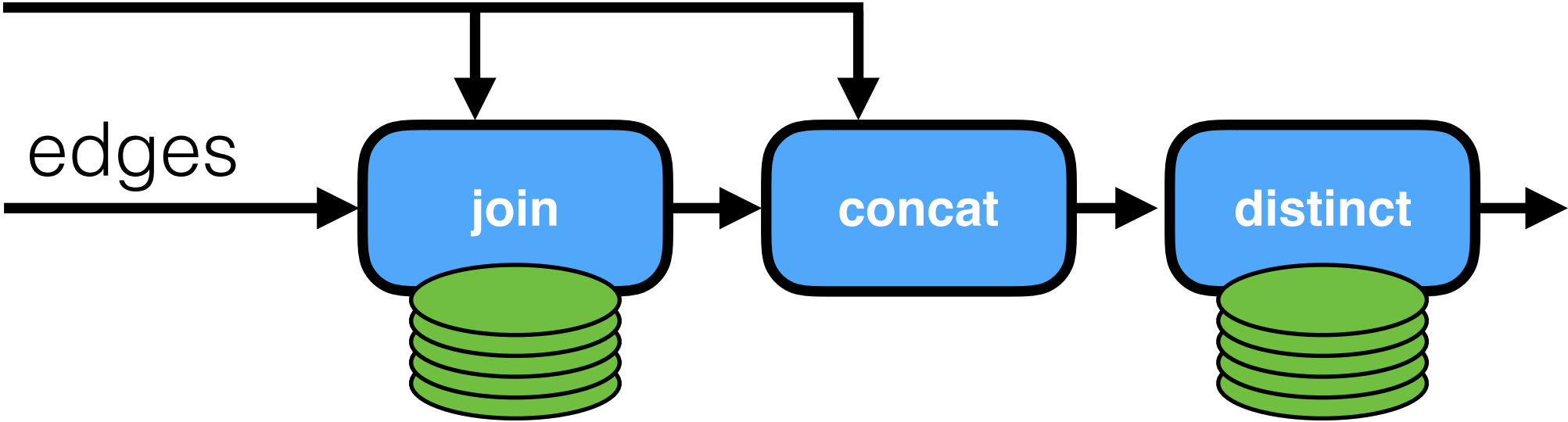
02us

Reach	cores	livejournal	orkut
GraphX	128	36s	48s
Socialite	128	52s	67s
Myria	128	5s	6s
BigDatalog	128	17s	20s
Differential	1	7s	15s
update	1	50us	62us

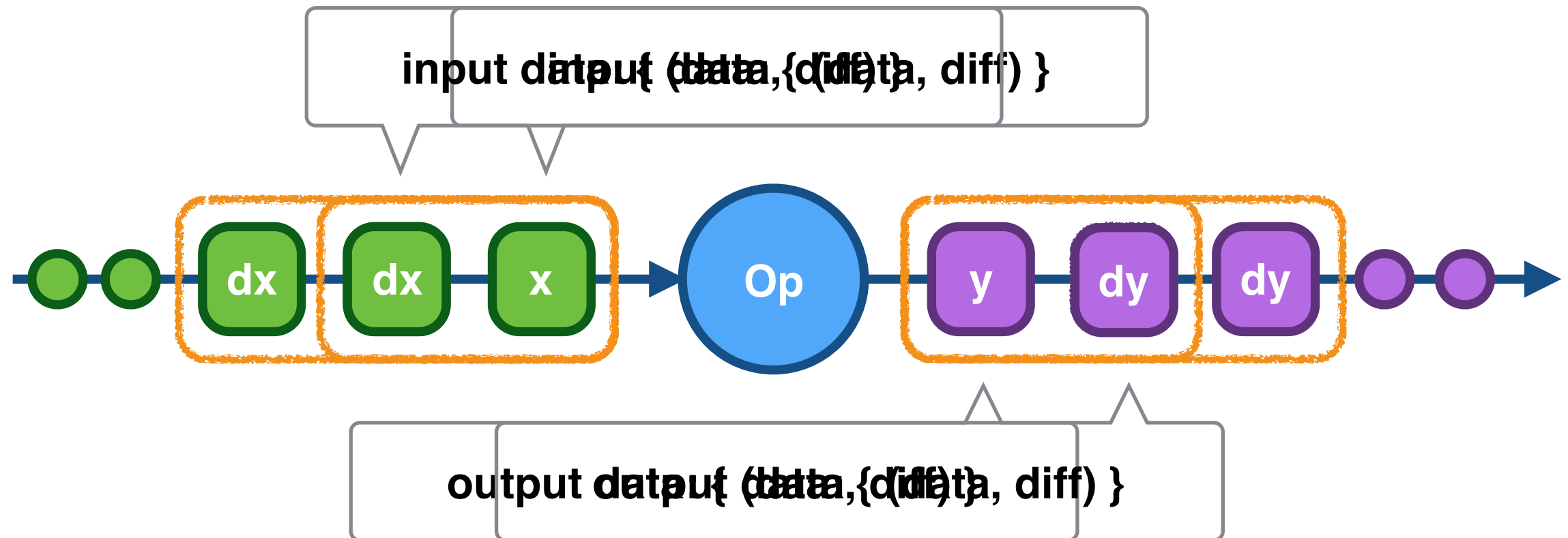
Stream<(Data, Time, isize)>

nodes

edges



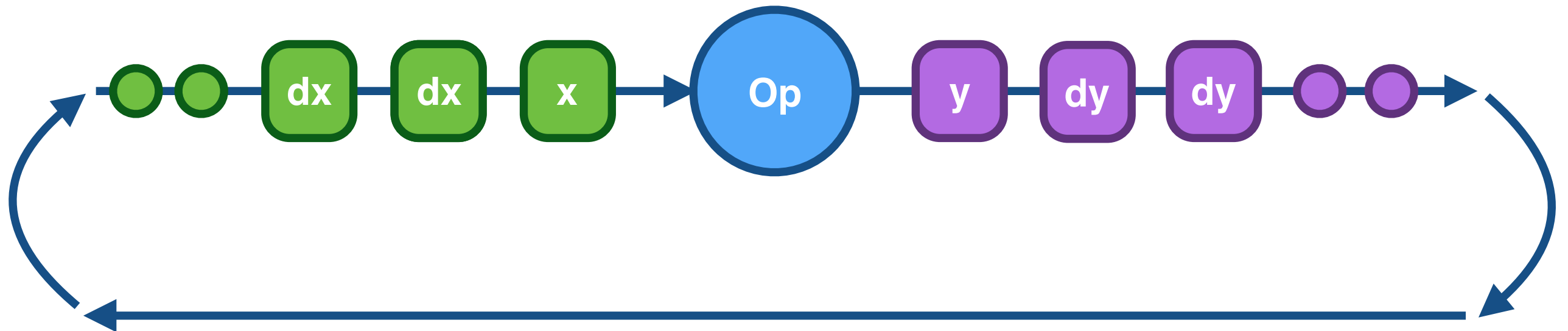
Incremental Dataflow



$$? = \text{Op} \left(x + dx \right) - y$$

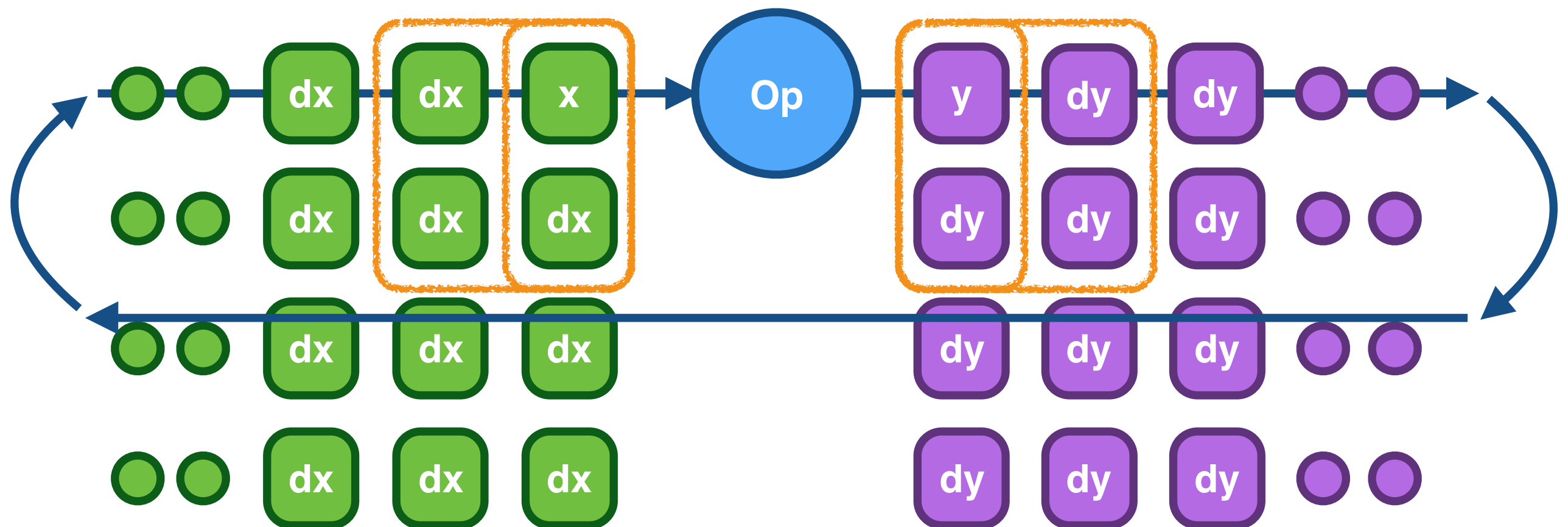
Iterative Dataflow

e.g. semi-naive bottom-up datalog

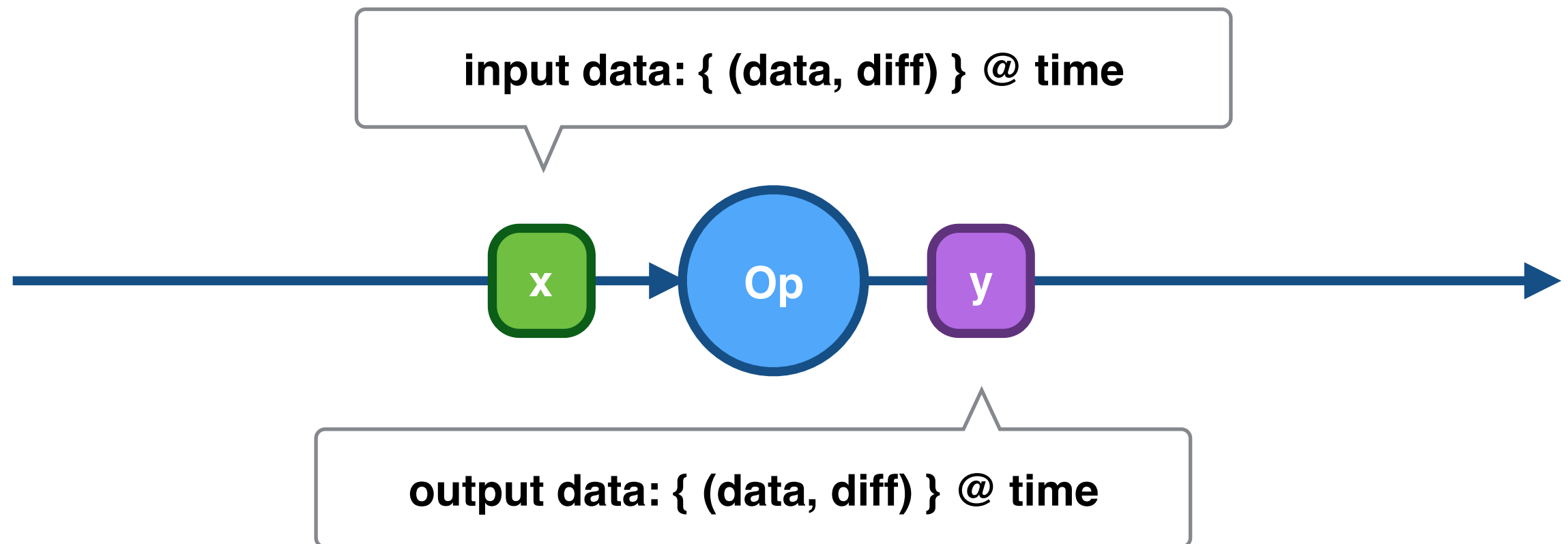


Differential Dataflow

(finally)



Differential Dataflow

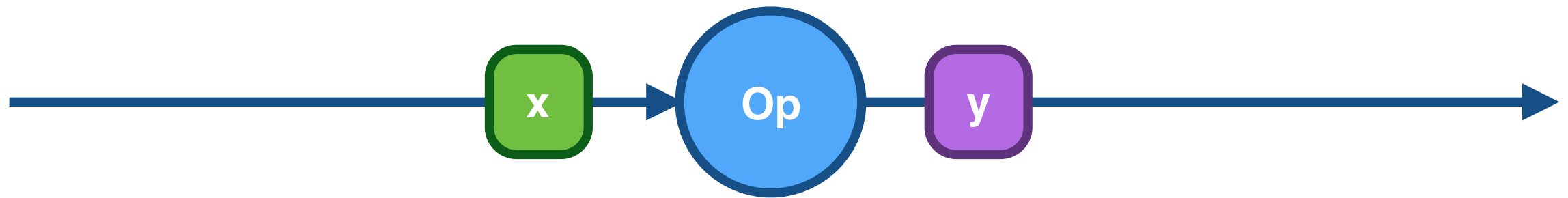


Important: times are only **partially** ordered

Differentiation on a discrete partial order

Differential Dataflow

with data-parallel operators



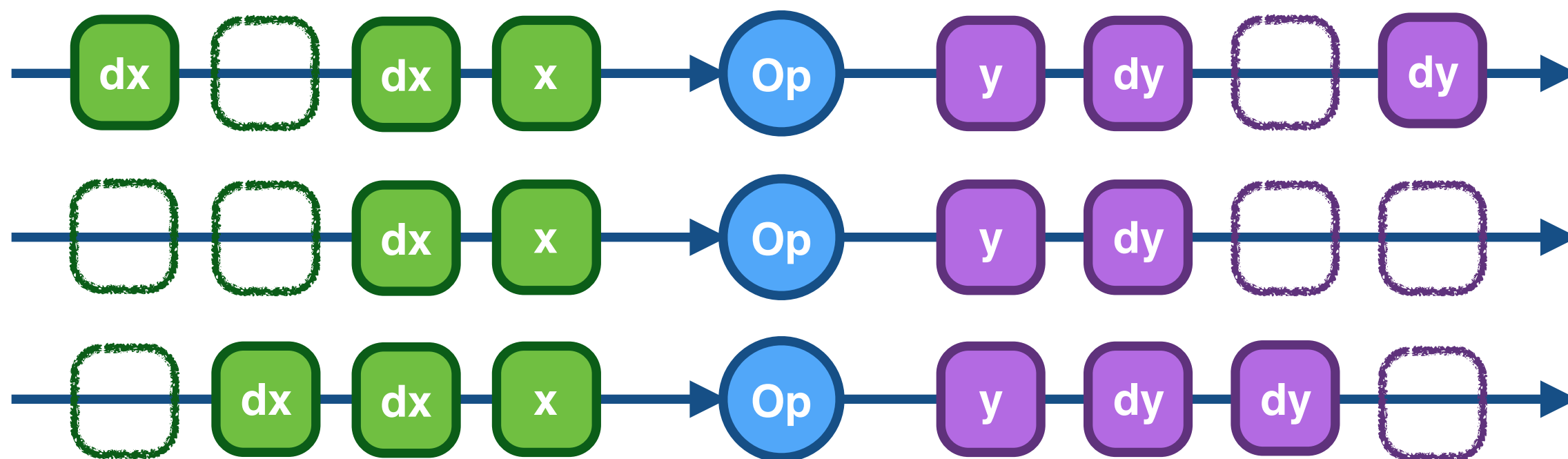
Differential Dataflow

with data-parallel operators



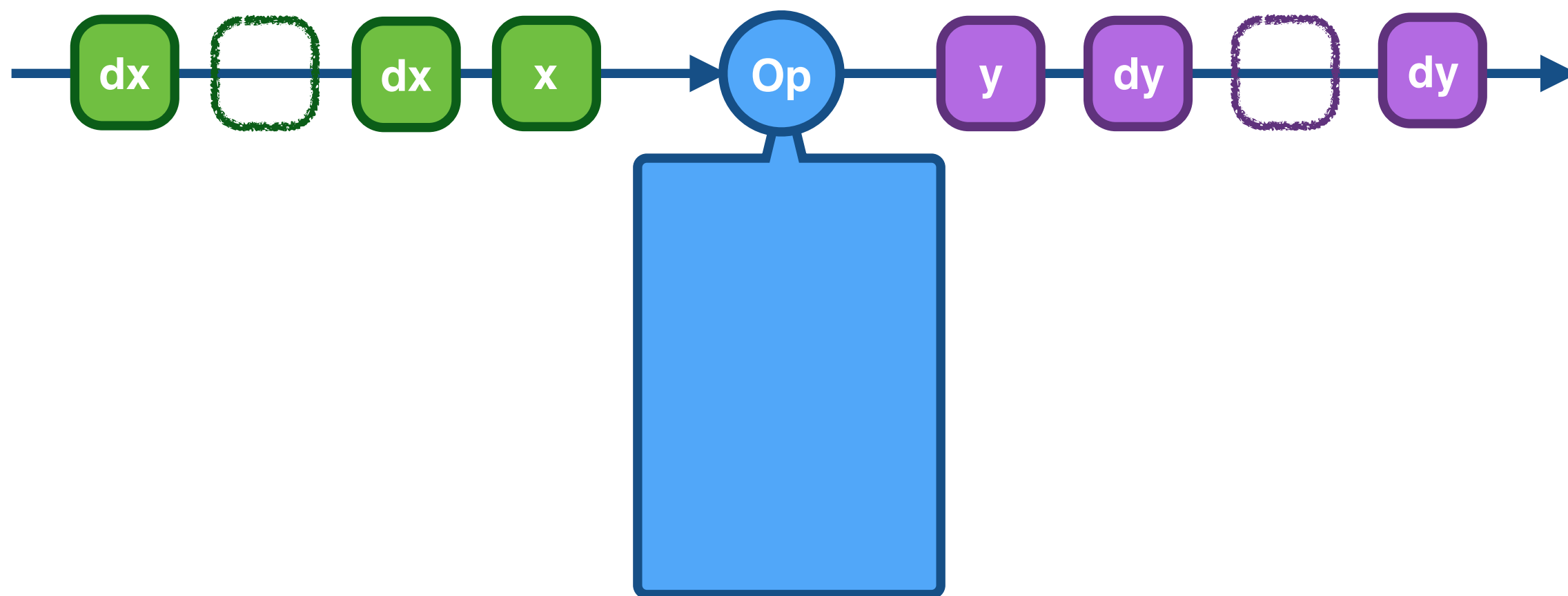
Differential Dataflow

with data-parallel operators



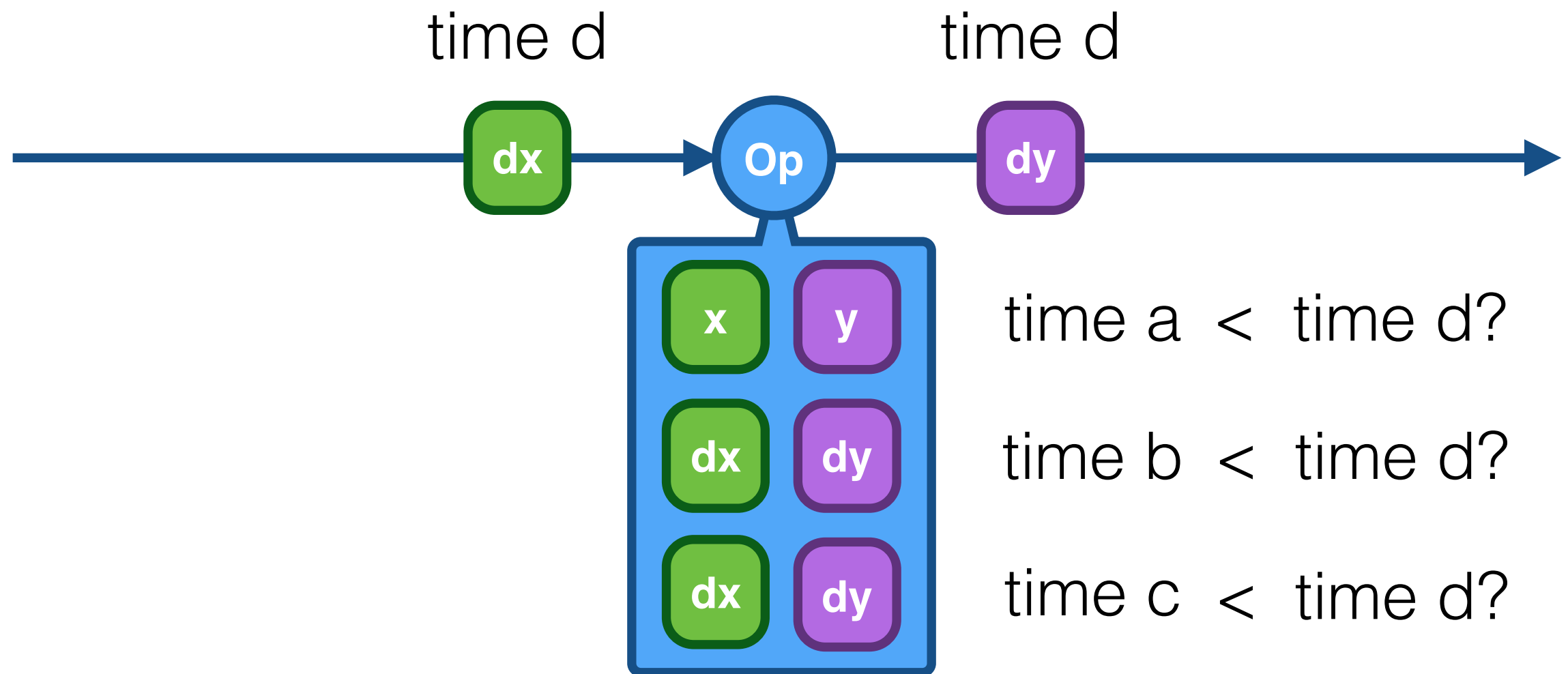
Differential Dataflow

with data-parallel operators



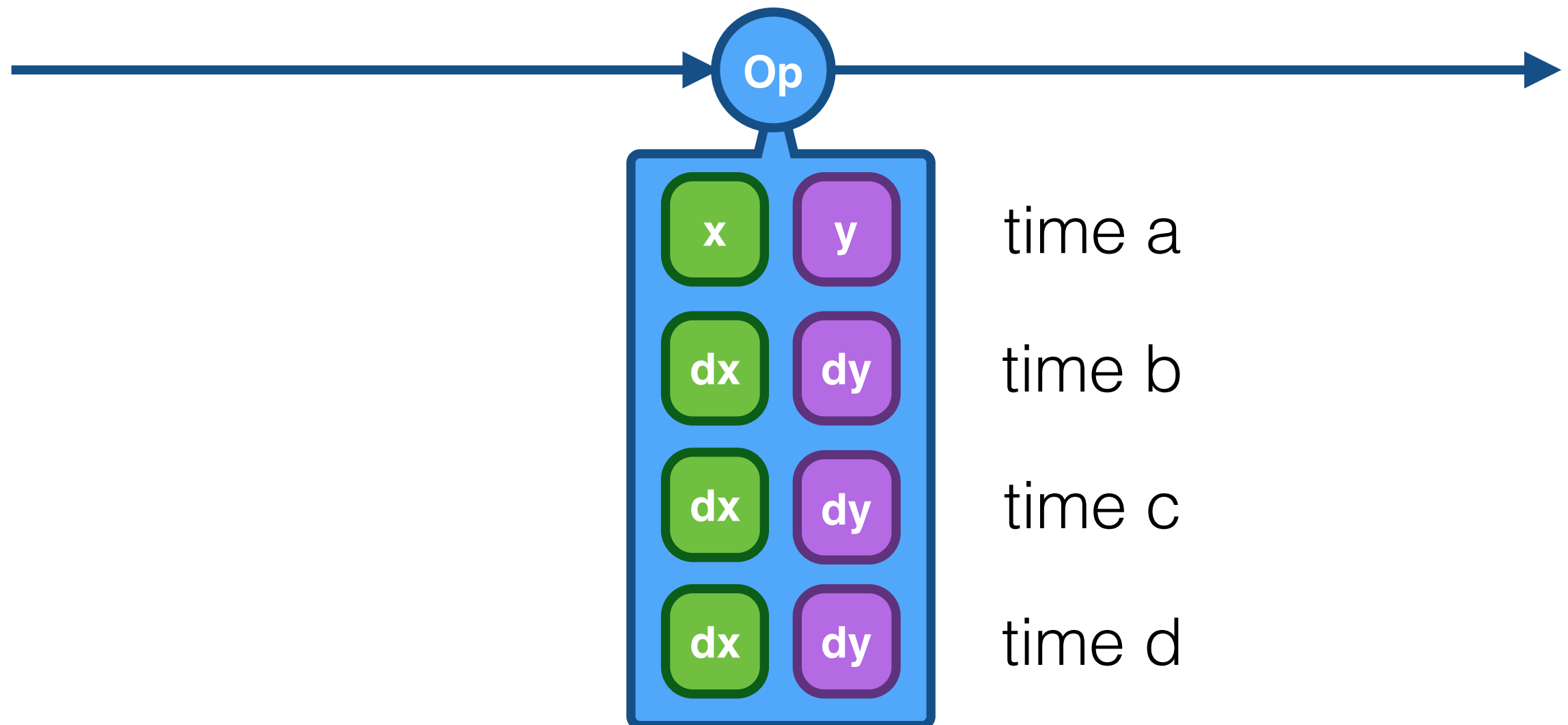
Differential Dataflow

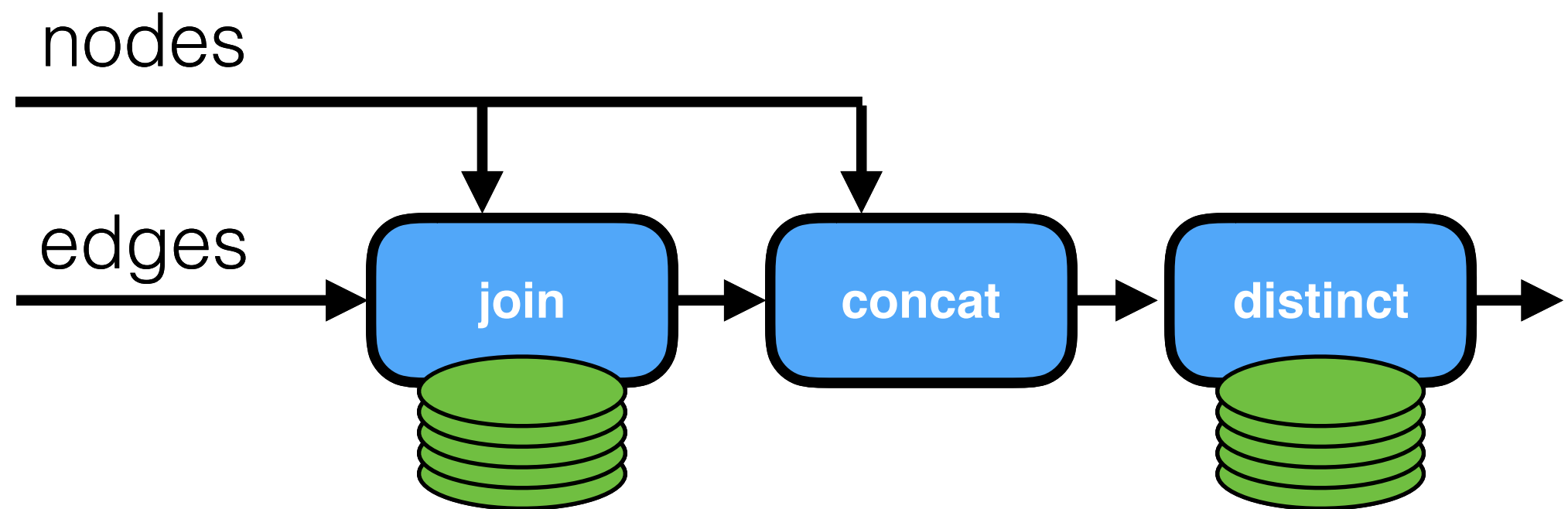
with data-parallel operators

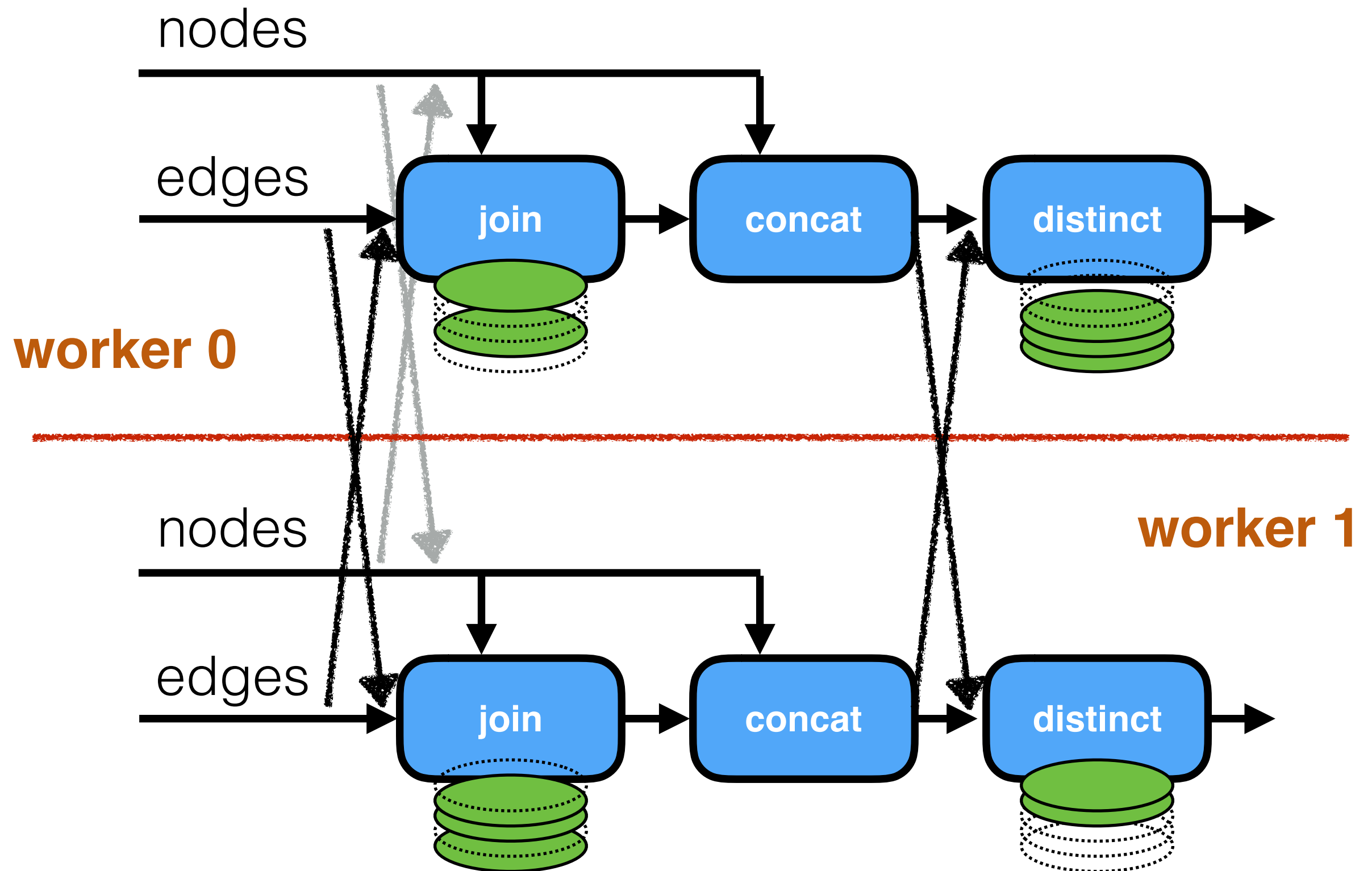


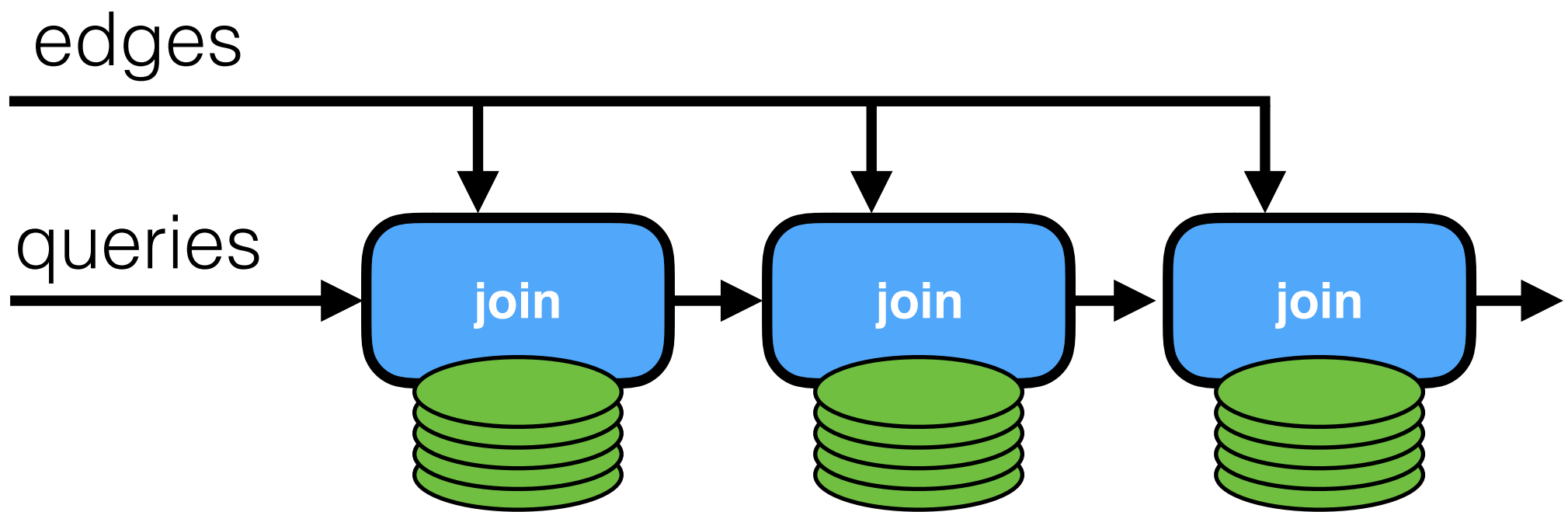
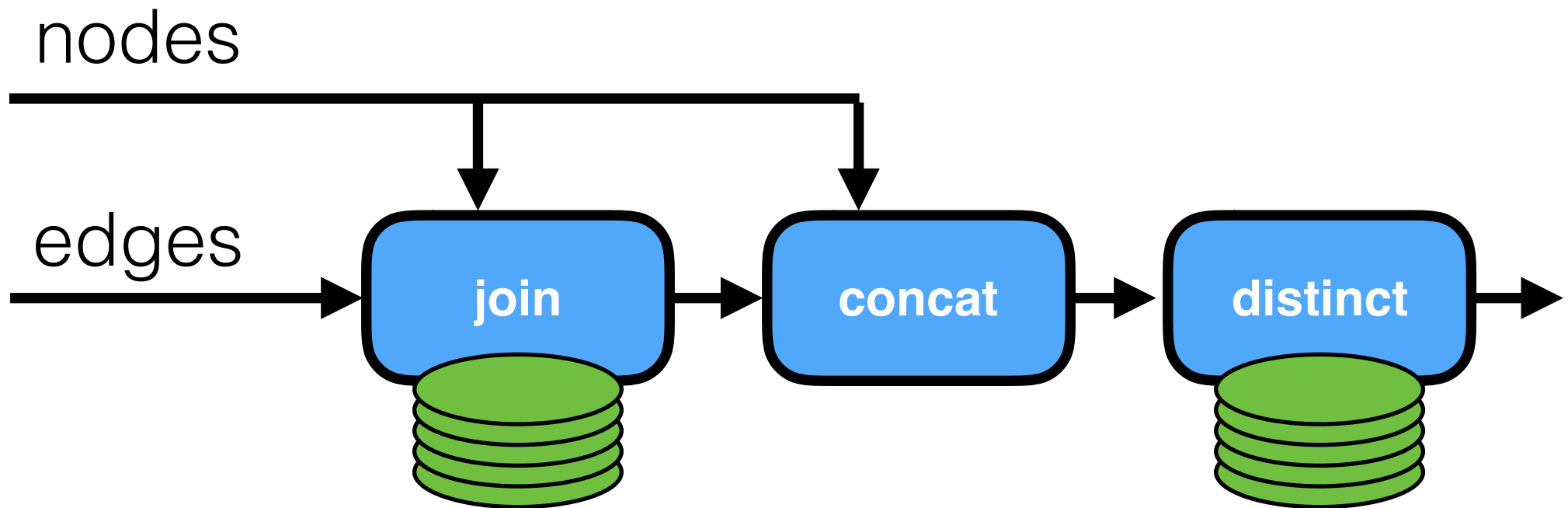
Differential Dataflow

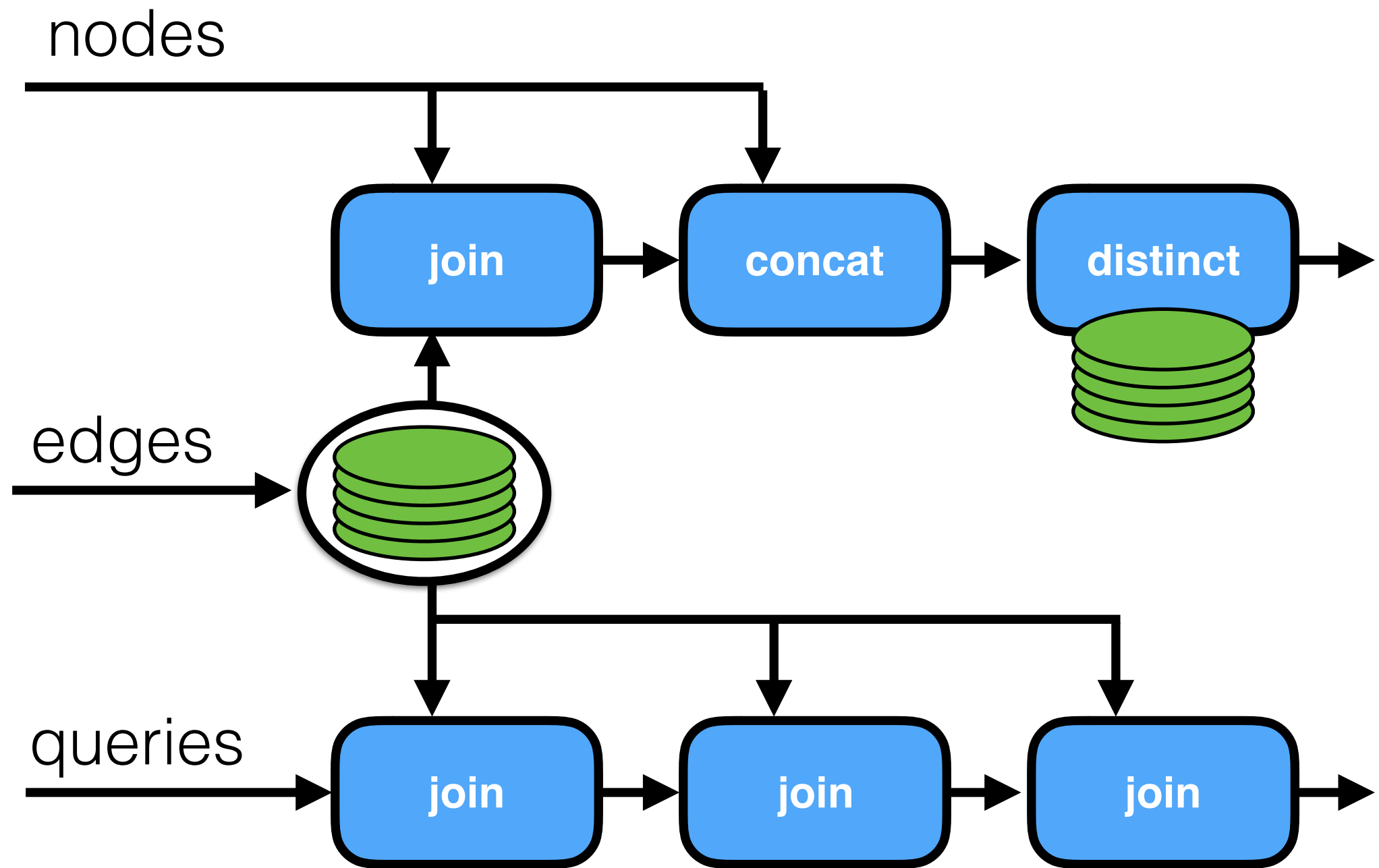
with data-parallel operators

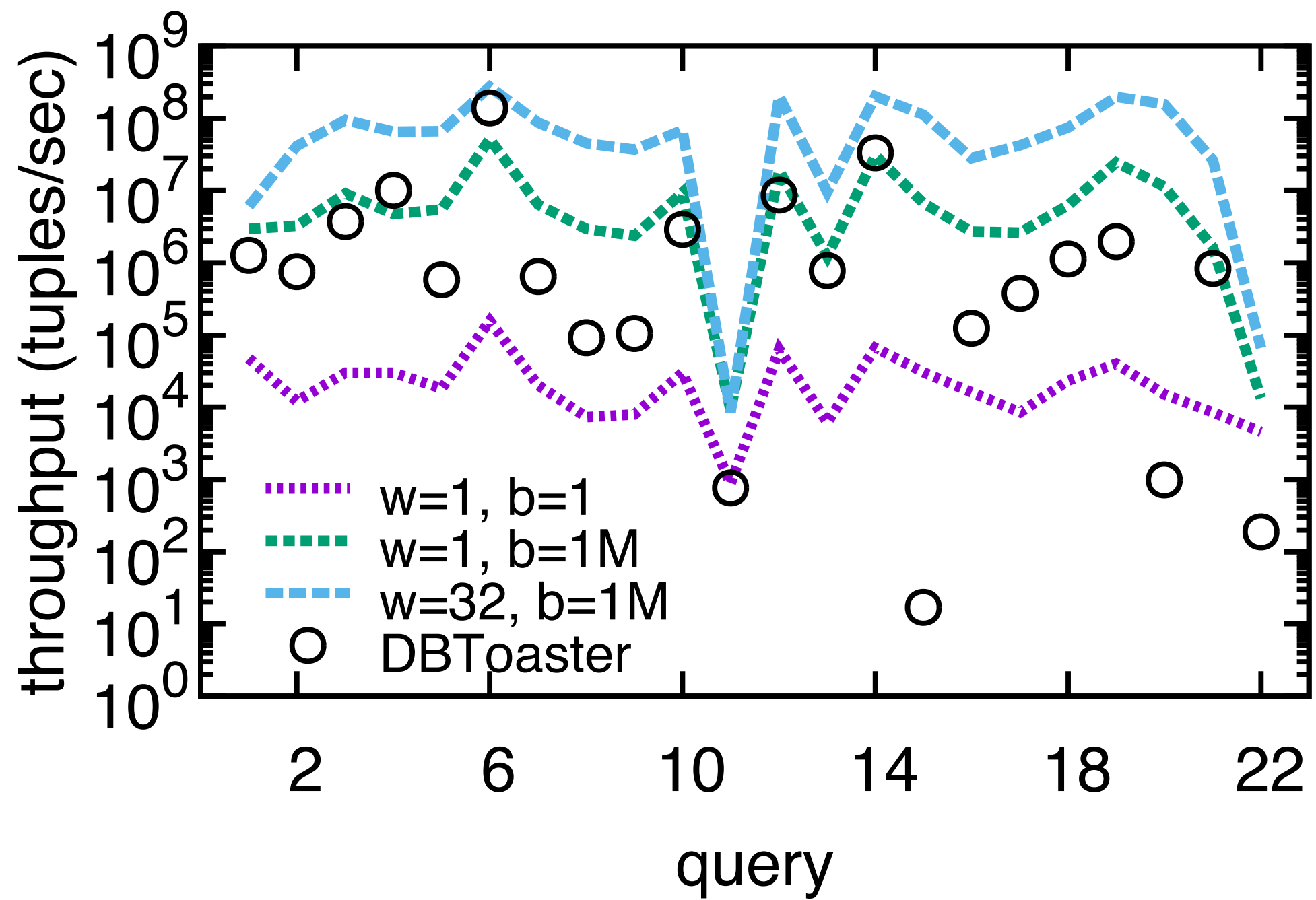


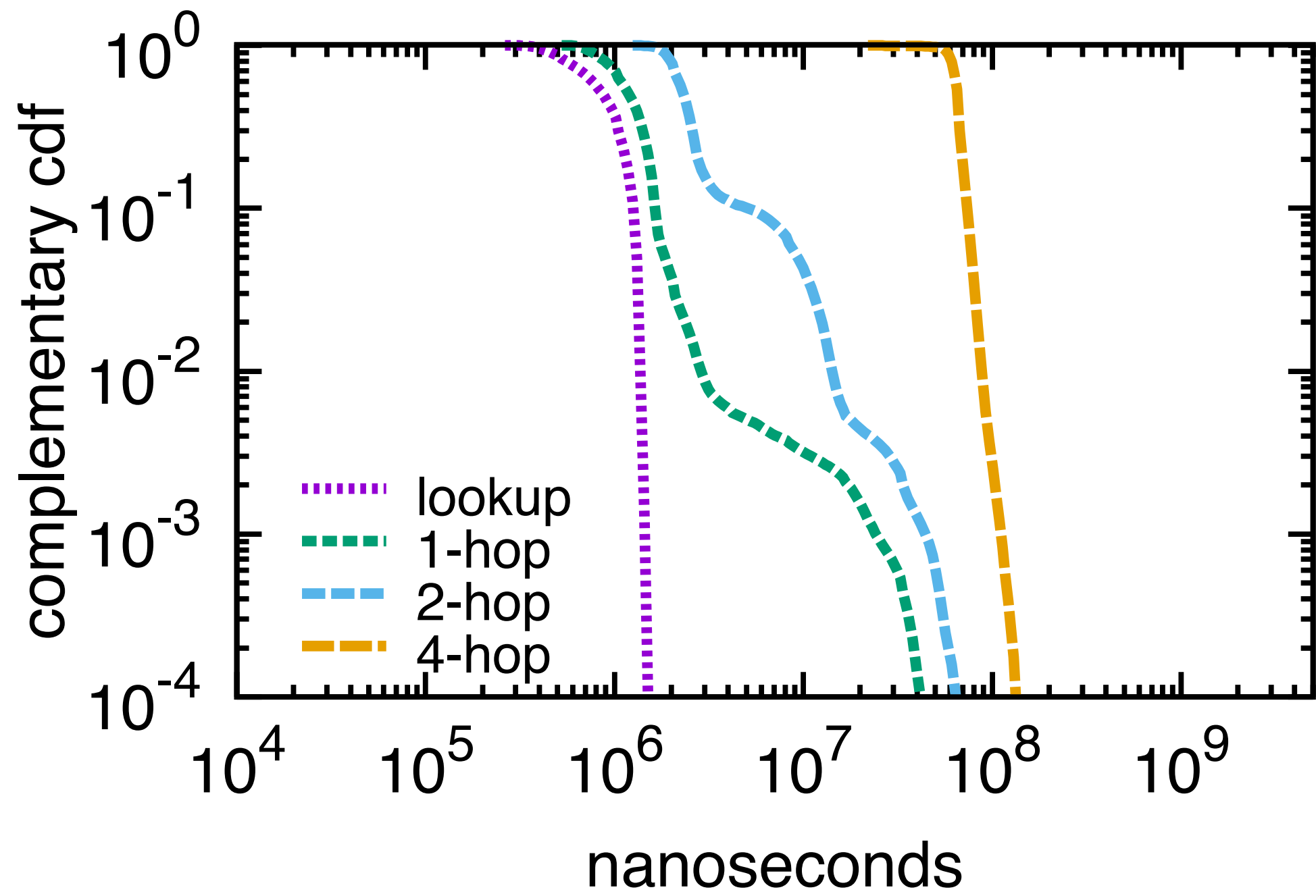












Dataflow	cores	httpd	psql	linux
Socialite	4	4h	OOM	OOM
Graspan	4	11.3m	143.8m	713.8m

[Graspan, ASPLOS 2017]

Dataflow	cores	httpd	psql	linux
Socialite	4	4h	OOM	OOM
Graspan	4	11.3m	143.8m	713.8m
Differential	1	10.9s	37.0s	76.8s
change (med)	1	22.0ms	185ms	1.11ms
change (max)	1	218ms	1.48s	8.13ms

PointsTo	cores	httpd	psql	linux
Socialite	4	>24h	OOM	OOM
Graspan	4	479.9m	353.1m	99.7m

[Graspan, ASPLOS 2017]

PointsTo	cores	httpd	psql	linux
Socialite	4	>24h	OOM	OOM
Graspan	4	479.9m	353.1m	99.7m
Differential	1	536.3s	362.0s	423.1s
Optimized	1	77.4s	75.9s	191.3s
-Sharing	1	91.9s	94.3s	401.7s

A quick taste:

Explaining outputs in modern data analytics
[VLDB2016]

Declarative Datalog debugging for mere mortals
[S. Köhler et al.]

“Describe provenance/lineage as Datalog queries”

Replace “Datalog” with Differential Dataflow:

1. More expressive (e.g. Map/Reduce/Aggr),
2. Incrementally maintained.

A quick taste:

Explaining outputs in modern data analytics
[VLDB2016]

Replace “Datalog” with Differential Dataflow:

1. More expressive (e.g. Map/Reduce/Aggr),
2. Incrementally maintained.

Explanations:

For a target subset of the observed output,
identify a subset of the input such that:
compute yields output containing the subset.

Differential Dataflow

current conclusions (ongoing)

Functional queries + timestamping = hugely powerful:

1. Automatic incrementalization (even w/loops).
2. Throughput scale-out w/o concurrency. (*c.f.* Calvin)
3. Shared state between multiple queries. (*c.f.* RDBMS)

Recovering RBDMS features w/o sacrificing usability.
We are sacrificing e.g. transactions. This is not OLTP.

Timely Dataflow

A data-parallel dataflow runtime

<http://github.com/frankmcsherry/timely-dataflow>

[based off of Naiad, SOSp 2013]

```
extern crate timely;

use timely::dataflow::operators::*;

fn main() {

    timely::example(|scope| {

        (0..10).to_stream(scope)
            .exchange(|&x| x)
            .inspect(|x| println!("{}", x));

    });

}
```

Your SPMD program.
You write this once.

Per worker logic:

ad.

computer.

```
extern crate timely;

use timely::dataflow::operators::*;

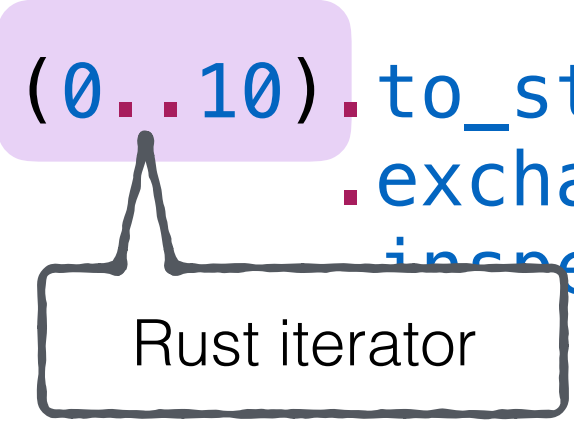
fn main() {

    timely::example(|scope| {

        (0..10).to_stream(scope)
            .exchange(|&x| x)
            .inspect(|x| println!("{}", x));

    });

}
```

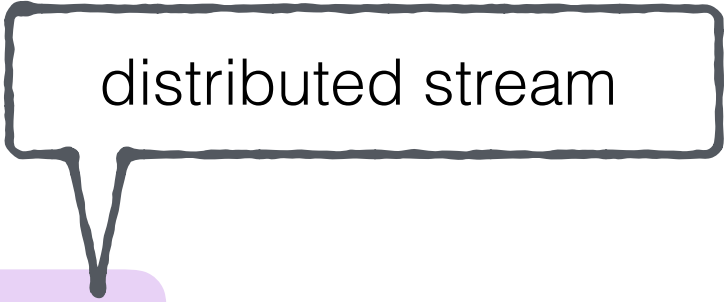


A diagram illustrating the flow of data in the provided Rust code. A light purple rounded rectangle highlights the expression `(0..10)` in the line `(0..10).to_stream(scope)`. A black line with an arrow points from this rectangle to a white rectangular box with a black border. Inside this box, the text "Rust iterator" is written in black. The arrow indicates that the `(0..10)` expression is a Rust iterator that is then converted into a stream by the `to_stream` method.

```
extern crate timely;

use timely::dataflow::operators::*;

fn main() {
    timely::example(|scope| {
        (0..10).to_stream(scope)
            .exchange(|&x| x)
            .inspect(|x| println!("{}", x));
    });
}
```

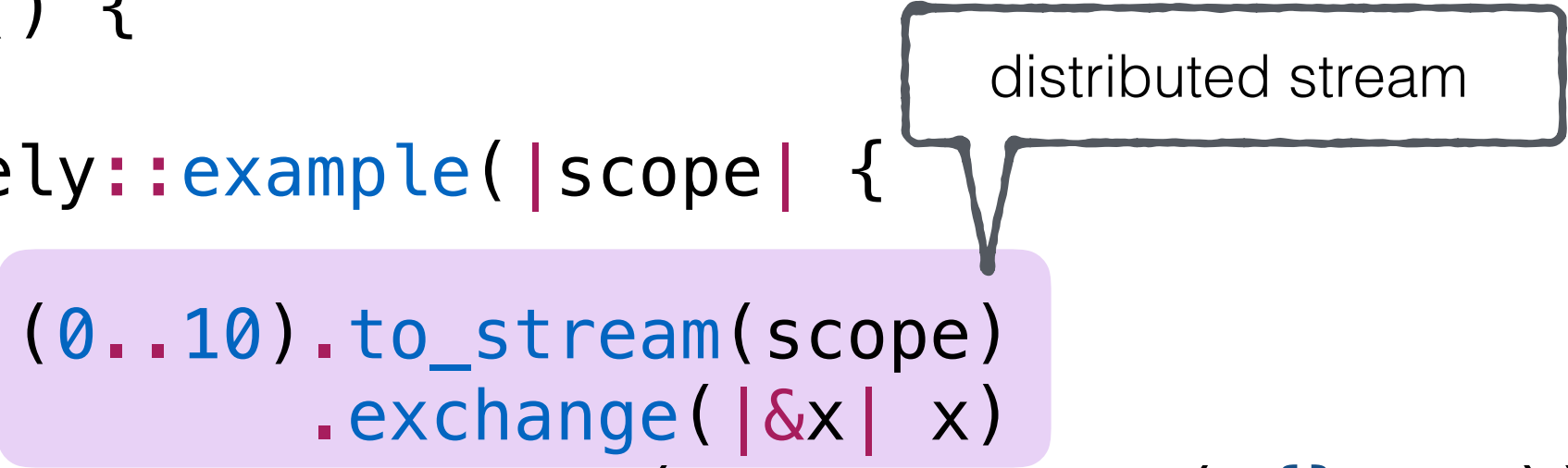


A callout box with a black border and a pointer to the `.to_stream(scope)` method call in the code. The text inside the box is "distributed stream".

```
extern crate timely;

use timely::dataflow::operators::*;

fn main() {
    timely::example(|scope| {
        (0..10).to_stream(scope)
            .exchange(|&x| x)
            .inspect(|x| println!("{}", x));
    });
}
```



A callout box with a pointer to the `example` function in the code, containing the text "distributed stream". A light purple rectangular block highlights the stream processing logic: `(0..10).to_stream(scope).exchange(|&x| x).inspect(|x| println!("{}", x));`.

```
extern crate timely;
```

```
use timely::dataflow::operators::*;
```

```
fn main() {
```

```
    timely::example(|scope| {
```

distributed stream

```
        (0..10).to_stream(scope)
            .exchange(|&x| x)
            .inspect(|x| println!("{}", x));
```

```
    });
```

```
}
```

```
fn main() {  
    timely::execute_from_args(std::env::args(), |worker| {
```

```
    }  
});  
}
```



```
fn main() {  
  timely::execute_from_args(std::env::args(), |worker| {
```

Adds a timestamp to each record

```
    let input = worker.new_input();  
    let probe = worker.dataflow(|scope| {  
      input.to_stream(scope)  
        .exchange(|&x| x)  
        .inspect(|x| println!("{}", x))  
        .probe()  
    });
```

Reports on possibility of timestamps

graph using a
new input and ending in a probe.

```
});
```

```
}
```

```

fn main() {
    timely::execute_from_args(std::env::args(), |worker| {

        let input = worker.new_input();
        let probe = worker.dataflow(|scope| {
            input.to_stream(scope)
                .exchange(|&x| x)
                .inspect(|x| println!("{}", x))
                .probe()
        });

        for round in 0..10 {
            input.send(round * round);
            input.advance_to(round + 1);
            while probe.less_than(input.time()) {
                worker.step();
            }
        }
    });
}

```

Drive the dataflow by supplying inputs and running until cleared.

“What times might **probe** still see?”
Information about distributed state.

```

fn main() {
    timely::execute_from_args(std::env::args(), |worker| {

        let input = worker.new_input();
        let probe = worker.dataflow(|scope| {
            input.to_stream(scope)
                .exchange(|&x| x)
                .inspect(|x| println!("{}", x))
                .probe()
        });

        for round in 0..10 {
            input.send(round * round);
            input.advance_to(round + 1);
            while probe.less_than(input.time()) {
                worker.step();
            }
        }
    });
}

```

```

        .iterate(|cycle| {
            .filter(|&x| x > 0)
            .map(|x| x - 1)
            .exchange(|&x| x)
        })

```

exchange benchmark

one thread:	1us
two threads:	2us
two processes:	40us

```

fn main() {
    timely::execute_from_args(std::env::args(), |worker| {

        let input = worker.new_input();
        let probe = worker.dataflow(|scope| {
            input.to_stream(scope)
                .probe()
        });

        for round in 0..10 {
            input.send(round * round);
            input.advance_to(round + 1);
            while probe.less_than(input.time() - 2) {
                worker.step();
            }
        }
    });
}

```

← -----

```

        .iterate(|cycle| {
            cycle.filter(|&x| x > 0)
                .map(|x| x - 1)
                .exchange(|&x| x)
        })

```

concurrent iterations

better utilizations!

```
// convenient!  
stream.map(logic);
```

Data movement requirements.
Could have been **Pipeline**.

behavior when scheduled.

```
// convenient!  
stream.map(logic);
```

Clarity on resource management:
Ownership drives collection/reuse.

```
// equivalent  
stream.unfold(
```

Destructuring ensures valid data maintained, reused

```
while let Some((time, data)) = input.next() {  
    let  
    for
```

It's basically just your code running.

```
    session.give(logic
```

Required to send

```
    drain(..)
```

Memory / control safety at compile time.
Fewer errors, more predictable behavior.

```
});
```

Abomonation

Serialization at memory bandwidth

<https://github.com/frankmcsherry/abomonation>

The name is from the first
reaction to the approach.

Be warned.


```
// serializes referenced typed data into bytes.
fn encode<T: Abomonation>(typed: &T, bytes: &mut Vec<u8>);

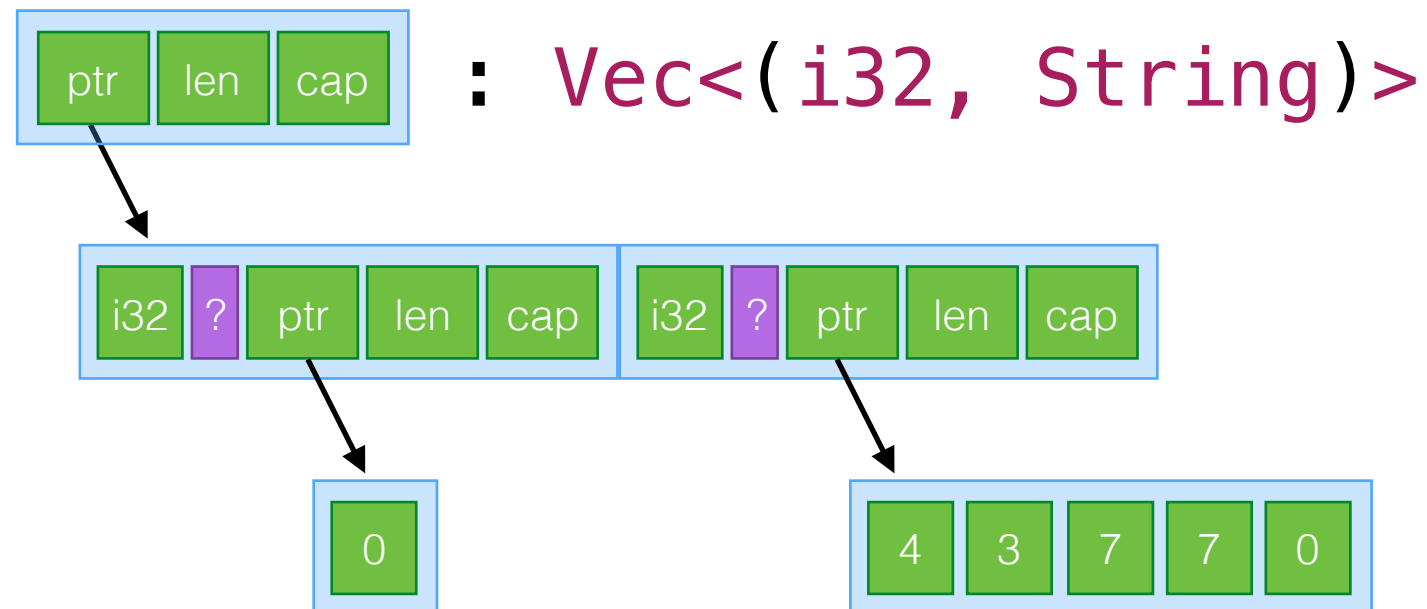
// presents a typed view of serialized data.
fn decode<T: Abomonation>(bytes: &mut [u8]) -> &T;
```

bytes:

What could this reference other than bytes?



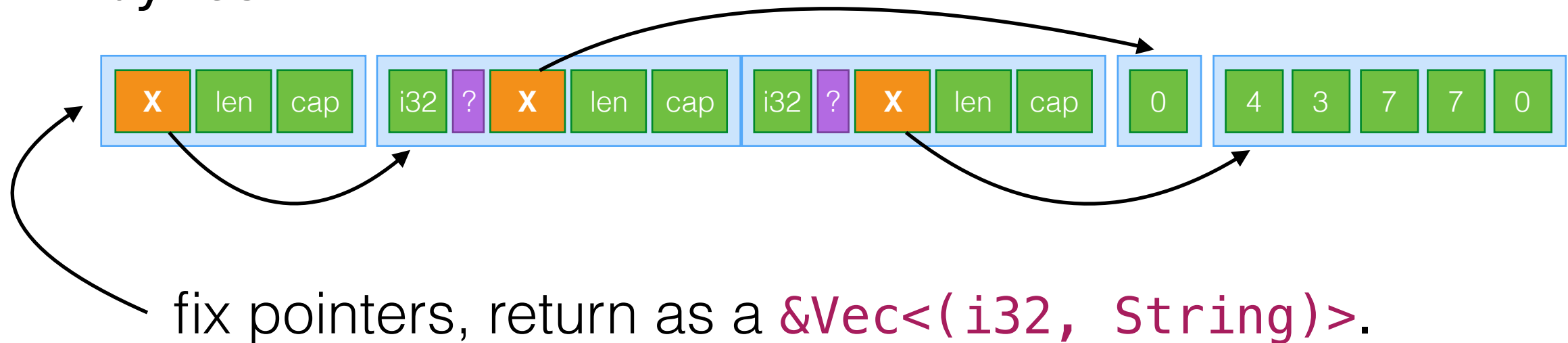
typed:



```
// serializes referenced typed data into bytes.
fn encode<T: Abomonation>(typed: &T, bytes: &mut Vec<u8>);

// presents a typed view of serialized data.
fn decode<T: Abomonation>(bytes: &mut [u8]) -> &T;
```

bytes:



1. totally ignores alignment (x64: ok, arm: bad).
2. memcpy of padding bytes is apparently UB.
3. assumes same architecture, binary.
4. blah blah blah.

```
impl Log {
  pub fn new() -> Log {
    Log {
      timestamp: 2837513946597,
      zone_id: 123456,
      zone_plan: ZonePlan::FREE,
      http: Http {
        protocol: HttpProtocol::HTTP11,
        status: 200,
```

```
unsafe_abominate!(Log : http, origin, server_ip, server_name, remote_ip, ray_id);
unsafe_abominate!(Http : content_type, user_agent, referer, request_uri);
unsafe_abominate!(Origin : ip, hostname);
```

```
      referer: "https://www.cloudflare.com".to_owned(),
      request_uri: "/cdn-cgi/trace".to_owned(),
    },
    origin: Origin {
      ip: "1.2.3.4".to_owned(),
      port: 8000,
      hostname: "www.example.com".to_owned(),
      protocol: OriginProtocol::HTTPS,
    },
    country: Country::US,
    cache_status: CacheStatus::Hit,
    server_ip: "192.168.1.1".to_owned(),
    server_name: "metal.cloudflare.com".to_owned(),
    remote_ip: "10.1.2.3".to_owned(),
    bytes_dlv: 123456,
    ray_id: "10c73629cce30078-LAX".to_owned(),
  }
}
```

```
1  #![feature(test)]
2  #[macro_use]
3  extern crate abomonation;
4  extern crate test;
5
6  use test::Bencher;
7  use abomonation::{Abomonation, encode, decode};
8
9  #[bench]
10 fn bench_populate(b: &mut Bencher) {
11     b.iter(|| {
12         Log::new()
13     });
14 }
```

bench_populate:

267 ns/iter (+/- 124)

```

9  #[bench]
10 fn bench_populate(b: &mut Bencher) {
11     b.iter(|| {
12         Log::new()
13     });
14 }
15
16 #[bench]
17 fn bench_serialize(b: &mut Bencher) {
18     let log = Log::new();
19     let mut bytes = vec![];
20     unsafe { encode(&log, &mut bytes); }
21     b.bytes = bytes.len() as u64;
22     b.iter(|| {
23         bytes.clear();
24         unsafe { encode(&log, &mut bytes); }
25         test::black_box(&bytes);
26     });
27 }

```

bench_populate: 267 ns/iter (+/- 124)

```

9  #[bench]
10 fn bench_populate(b: &mut Bencher) {
11     b.iter(|| {
12         Log::new()
13     });
14 }
15
16 #[bench]
17 fn bench_serialize(b: &mut Bencher) {
18     let log = Log::new();
19     let mut bytes = vec![];
20     unsafe { encode(&log, &mut bytes); }
21     b.bytes = bytes.len() as u64;
22     b.iter(|| {
23         bytes.clear();
24         unsafe { encode(&log, &mut bytes); }
25         test::black_box(&bytes);
26     });
27 }

```

bench_populate:	267 ns/iter (+/- 124)	
bench_serialize:	54 ns/iter (+/- 4)	= 10148 MB/s

```

21     b.bytes = bytes.len() as u64;
22     b.iter(|| {
23         bytes.clear();
24         unsafe { encode(&log, &mut bytes); }
25         test::black_box(&bytes);
26     });
27 }
28
29 #[bench]
30 fn bench_deserialize(b: &mut Bencher) {
31     let log = Log::new();
32     let mut bytes = vec![];
33     unsafe { encode(&log, &mut bytes); }
34     b.bytes = bytes.len() as u64;
35     b.iter(|| {
36         test::black_box(unsafe { decode::<Log>(&mut bytes) });
37     });
38 }

```

bench_populate:	267 ns/iter (+/- 124)
bench_serialize:	54 ns/iter (+/- 4) = 10148 MB/s

```

21     b.bytes = bytes.len() as u64;
22     b.iter(|| {
23         bytes.clear();
24         unsafe { encode(&log, &mut bytes); }
25         test::black_box(&bytes);
26     });
27 }
28
29 #[bench]
30 fn bench_deserialize(b: &mut Bencher) {
31     let log = Log::new();
32     let mut bytes = vec![];
33     unsafe { encode(&log, &mut bytes); }
34     b.bytes = bytes.len() as u64;
35     b.iter(|| {
36         test::black_box(unsafe { decode::<Log>(&mut bytes) });
37     });
38 }

```

bench_populate:	267 ns/iter (+/- 124)	
bench_serialize:	54 ns/iter (+/- 4)	= 10148 MB/s
bench_deserialize:	8 ns/iter (+/- 1)	= 68500 MB/s


```

21     b.bytes = bytes.len() as u64;
22     b.iter(|| {
23         bytes.clear();
24         unsafe { encode(&log, &mut bytes); }
25         test::black_box(&bytes);
26     });
27 }
28
29 #[bench]
30 fn bench_deserialize(b: &mut Bencher) {
31     let log = Log::new();
32     let mut bytes = vec![];
33     unsafe { encode(&log, &mut bytes); }
34     b.bytes = bytes.len() as u64;
35     b.iter(|| {
36         test::black_box(unsafe { decode::<Log>(&mut bytes) });
37     });
38 }

```

It's even faster than this

This is really fast

bench_populate:	267 ns/iter (+/- 124)	
bench_serialize:	54 ns/iter (+/- 4)	= 10148 MB/s
bench_deserialize:	8 ns/iter (+/- 1)	= 68500 MB/s
bench_deserialize_test:	112 ns/iter (+/- 25)	= 4892 MB/s

Projects

<http://github.com/frankmcsherry/timely-dataflow>

<http://github.com/frankmcsherry/differential-dataflow>

<http://github.com/frankmcsherry/abomonation>

Projects

<http://github.com/frankmcsherry/timely-dataflow>

<http://github.com/frankmcsherry/differential-dataflow>

<http://github.com/frankmcsherry/abomonation>

Moar projects

<http://github.com/frankmcsherry/recycler>

<http://github.com/frankmcsherry/columnar>

Projects

<http://github.com/frankmcsherry/timely-dataflow>

<http://github.com/frankmcsherry/differential-dataflow>

<http://github.com/frankmcsherry/abomonation>

Moar projects

<http://github.com/frankmcsherry/recycler>

<http://github.com/frankmcsherry/columnar>

Scandalous blog:

<http://github.com/frankmcsherry/blog>

Projects

<http://github.com/frankmcsherry/timely-dataflow>

<http://github.com/frankmcsherry/differential-dataflow>

<http://github.com/frankmcsherry/abomonation>

Moar projects

<http://github.com/frankmcsherry/recycler>

<http://github.com/frankmcsherry/columnar>

Scandalous blog:

<http://github.com/frankmcsherry/blog>

Excited to have more eyeballs, help, criticism. stars...