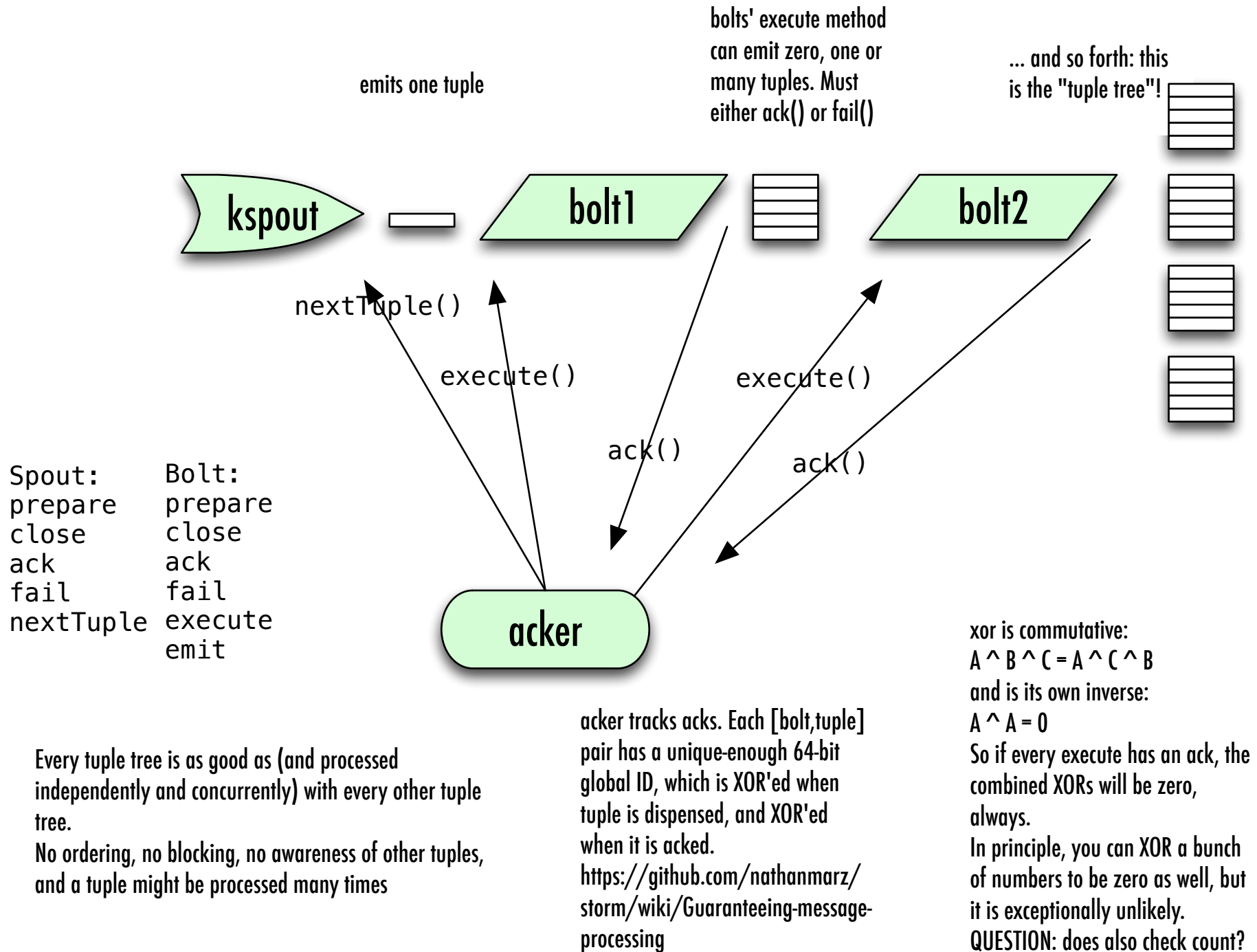
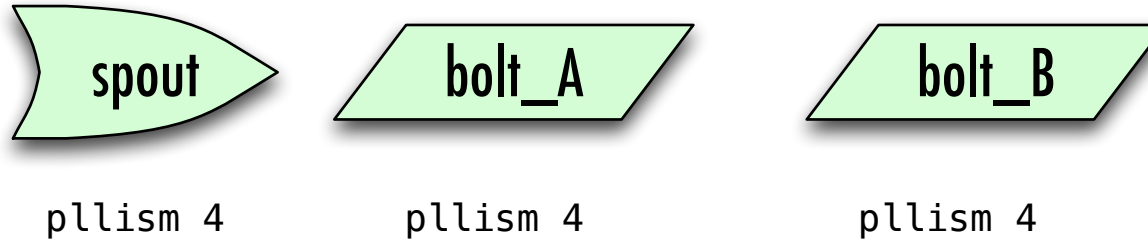


# Storm

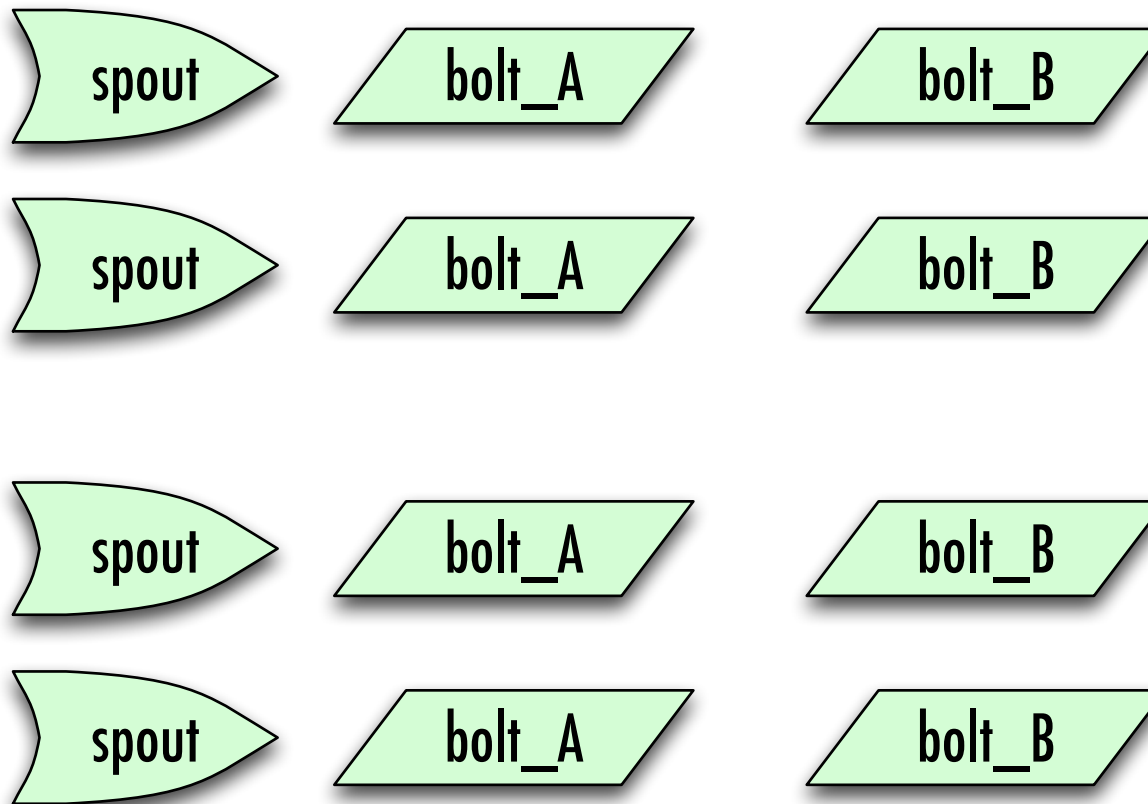


## Topology



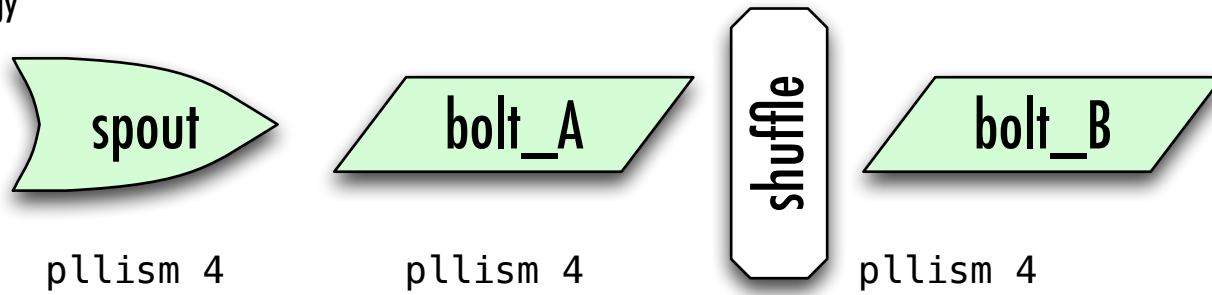
## Physical

These will send directly from spouts to A's to B's



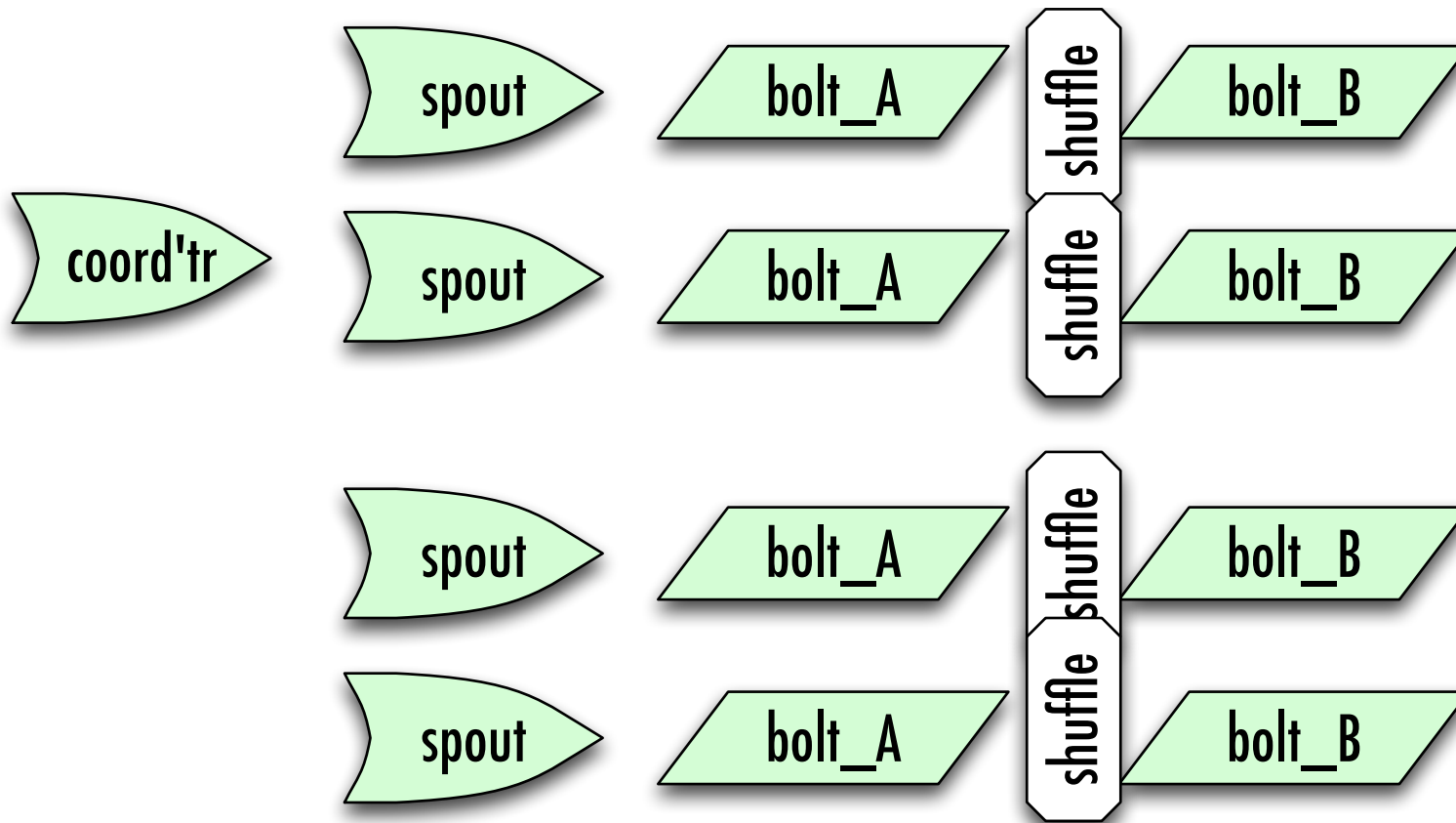


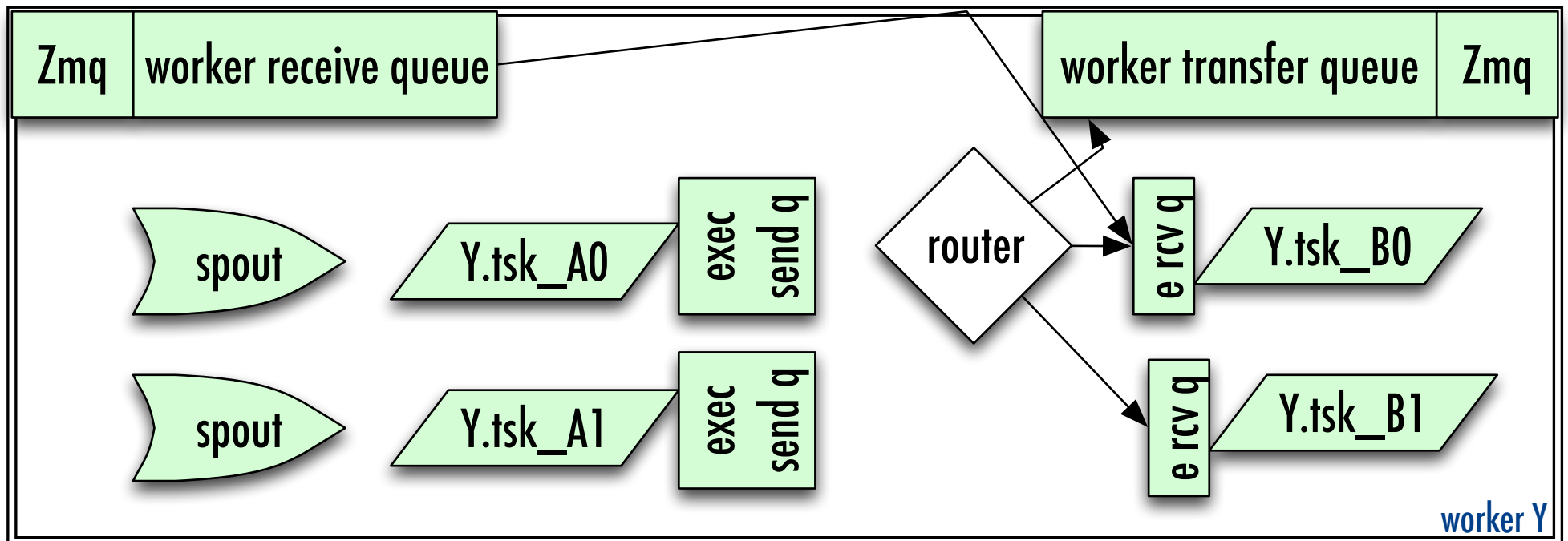
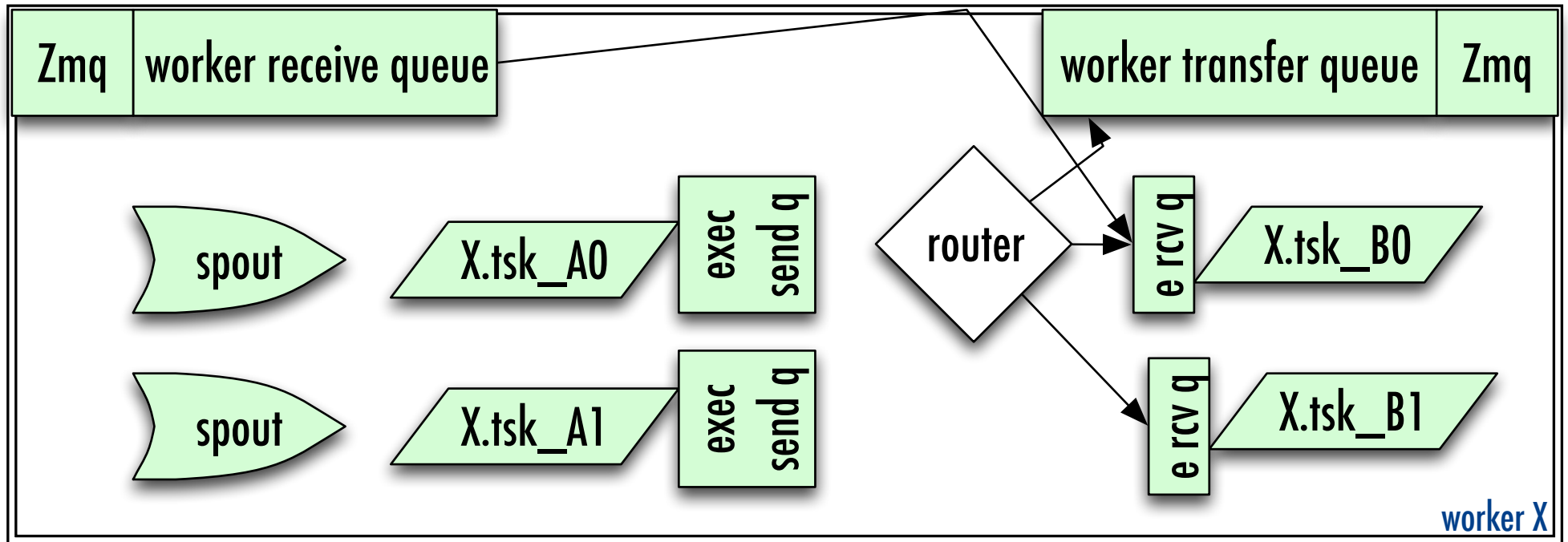
## Topology

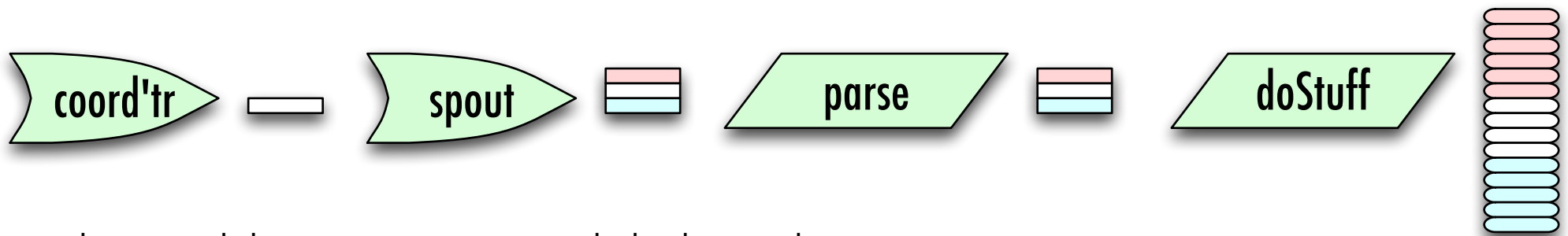


## Physical

These will need the send end receive buffers







coordinator is secretly the spout  
at trident batch delay period,  
will emit a transaction tuple  
it has a serially incrementing  
transaction ID, kept forever  
even across restarts.  
(We're only going to talk about  
Opaque Transactional topologies btw)

spout emits batches; these succeed  
or fail as a whole. (That's because  
they're part of the Storm tuple  
tree of the coordinator's seed  
tuple).

tuples in batches are freely processed in parallel and  
asynchronously unless there are barriers (eg the state after the  
group by)

TODO: where is the blocking behavior for state code

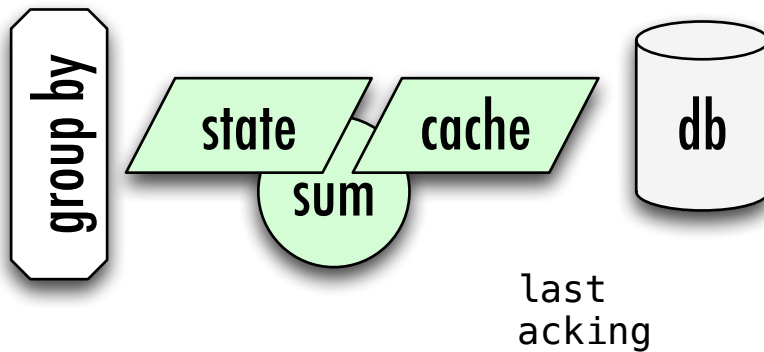
Trident:

- \* exactly once processing
- \* transactional db support
- \* layer on the flow DSL
- \* some primitive aggregators

non-transactional: batching behavior only

transactional: exactly once; batches are always  
processed in whole

opaque transactional: all records are processed, but  
might not be in same batches



the state doesn't ask the cache to fetch until it has a whole batches' worth of records to hand over. This is trident logic, not storm.

Those "agggregables" are reduced into rolled-up aggregates. So you might have 2500 inbound records that result in 900 distinct aggregates. (If you had eight aggregables [A, A, C, A, B, D, B, A] you would get four partial aggregates {A: 4, B: 2, C: 1, D: 1}).

It's clever about doing partial aggregates ("algebraic" reducers).

It looks in the cache for the old total count. Anything that isn't there it fetches from the database. This lets you do efficient batch requests, a huge scalability boon.

Once the cache is fresh, it determines the next aggregated value and writes it to the cache and to the DB, then ack()s the batch (all the tuples in the batch, really).

If a batch had 900 aggregates, and it had prior counts for 250 of them, then it will read 650 records and write 900. It always does a put for every new observed count.

¡Note!: The database writes do *\*not\** have to be transactional. It's the whole thing – the whole batch, end-to-end – that has to have transactional integrity, not just the DB code.

Let's say for transaction ID 69 the old aggregated values were {A: 24, B: 12, C: 2, D: 1}.

It stores

{A: [24, 20, 69], B: [12, 10, 69], C: [2, 1, 69], D: [1, 0, 69]}

TODD: verify order of that list

If I am processing batch

Since this is a State, you have contractual obligation from Trident to be processed until and unless batch 68 has succeeded.

So when I go to read from the DB, I will usually see something like

{A: [20, ??, 68], B: [10, ??, 68], C: [1, ??, 68]}

I might instead however see

{A: [??, 20, 69], B: [??, 10, 69], C: [??, 1, 69], D: [??, 0, 69]}

This means another attempt has been here: maybe it succeeded maybe I am the one who is succeeding but slow. In any case the new (first slot) values for this state, but I do know that I have values saved from batch 68. I just use those, and clobber the existing correct counts.

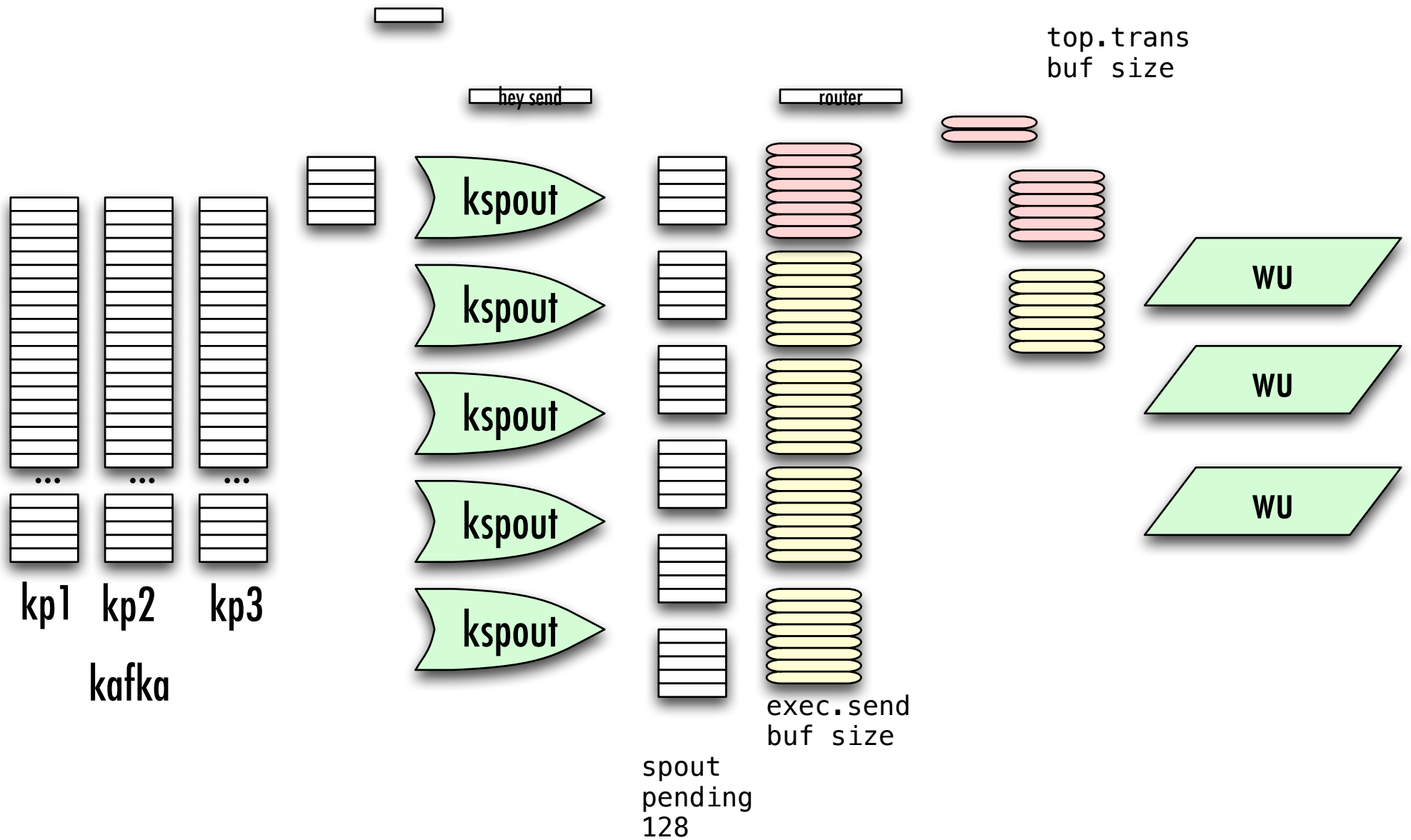
A:20, B: 10, C: 1, D: 0}, and new

ident that batch 69 will \*not\*

ike

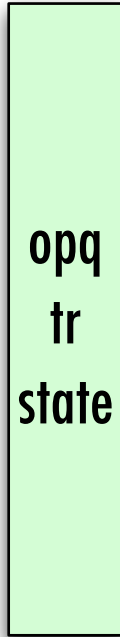
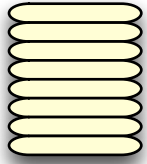
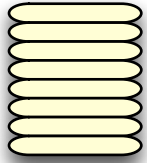
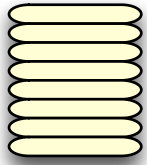
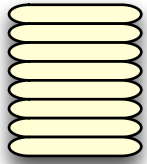
but was slow; maybe it failed;  
, I don't know whether to trust  
can trust the prior (second slot)  
isting values with my new,



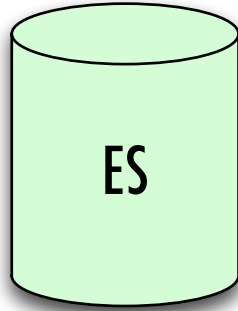


`max_fetch_size_bytes: 125k`  
(about 100 records)

`set spout||ism = kafka partitions`



opq  
tr  
state

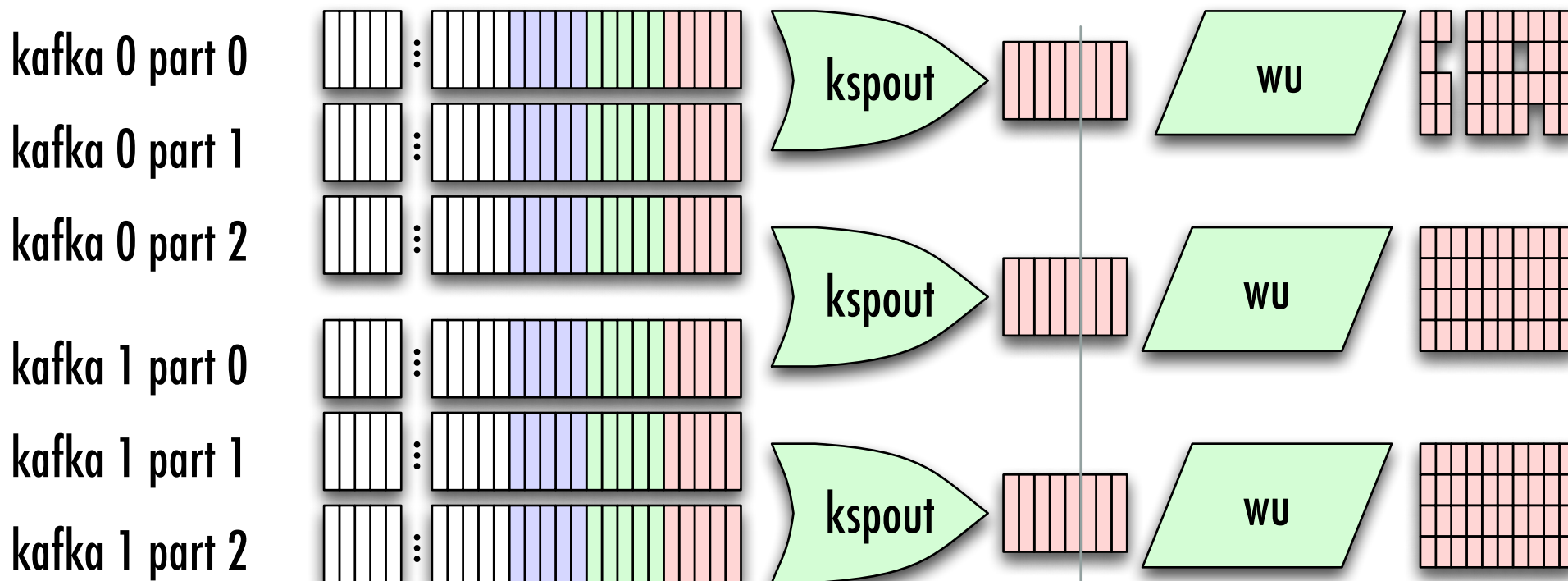


ES

exec.send  
buf size  
top.trans  
buf size

# Kafka

batch 1:0



must have more spouts than downstream execs if no shuffle

\* must have as many kfk partitions as spouts to keep them all fed

\* tasks an even multiple of executors;  
executors an even multiple of workers; partitions  
an even multiple of spouts; workers an even  
multiple of machines