

Programming Support for an Integrated Multi-Party Computation and MapReduce Infrastructure

Nikolaj Volgushev, Andrei Lapets, Azer Bestavros



Programming Support for an Integrated Multi-Party Computation and MapReduce Infrastructure



Programming Support for an Integrated **Multi-Party Computation** and **MapReduce** Infrastructure



Programming Support for an Integrated **Multi-Party Computation** and MapReduce Infrastructure



What's multi-party computation (MPC)?

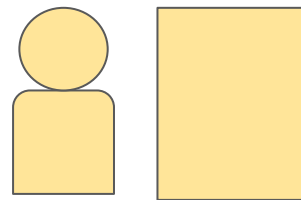
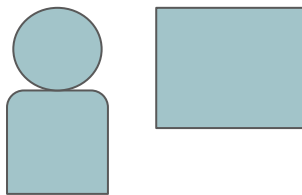
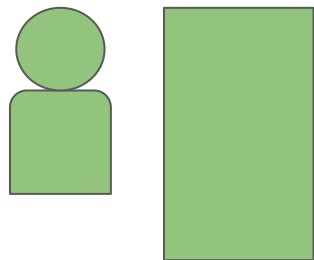
Given multiple parties p_1, p_2, \dots, p_n with private inputs x_1, x_2, \dots, x_n

Need to compute $f(x_1, x_2, \dots, x_n)$

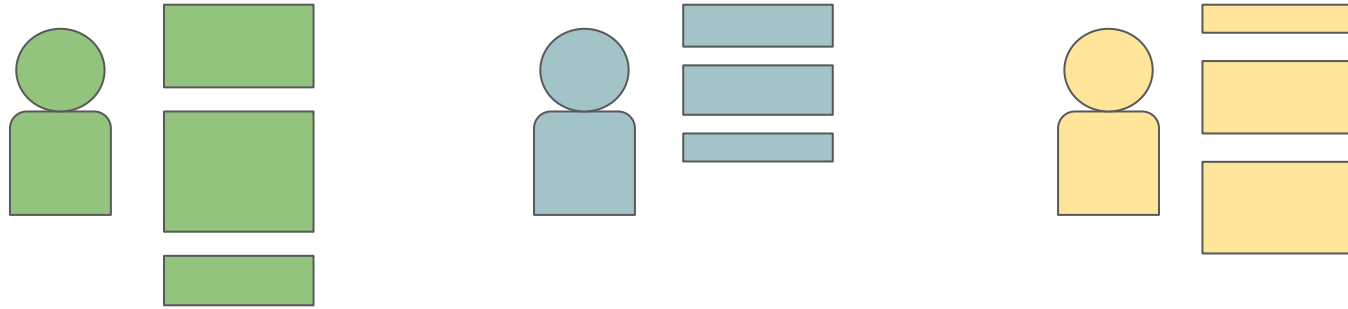
Without revealing more than the outputs of f

Sounds a bit like a magic trick.

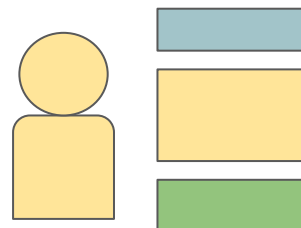
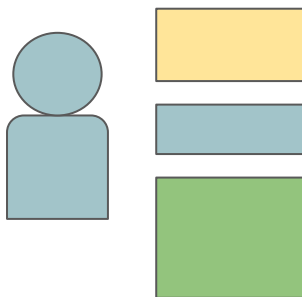
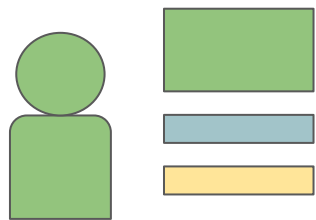
Quick example: the sum of secrets



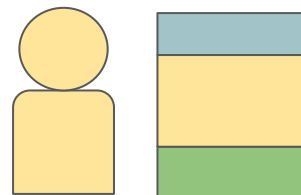
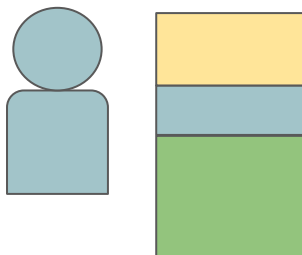
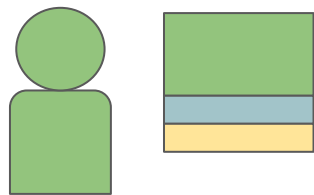
Players split their secrets into “shares”



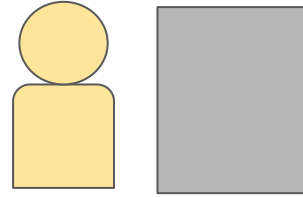
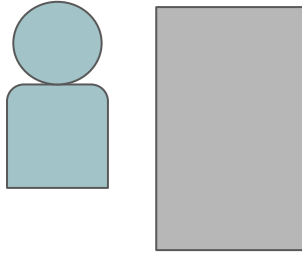
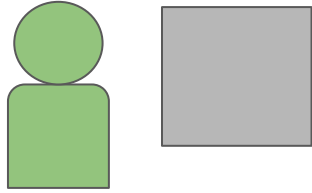
All players exchange shares



Each player computes the sum of their shares



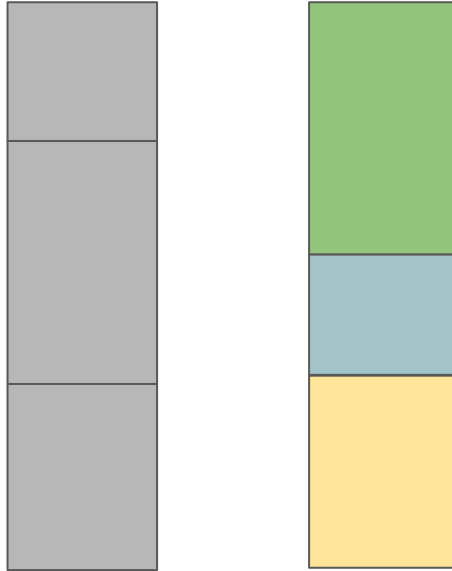
Each player now holds a share of the sum of secrets!



Exchange and recombine result shares



Lo and behold



MPC is great!

Let's us run computations, while preserving privacy of the inputs.

A mathematical turtle shell for our data!



MPC is great (in theory)!

Let's us run computations, while preserving privacy of the inputs.

A mathematical turtle shell for our data!

Practical MPC frameworks exist but they are **slow** (not unlike turtles)

The learning curve is steep (trust me!)



Programming Support for an Integrated Multi-Party Computation and **MapReduce** Infrastructure



MapReduce: fast, like a hare!

Programming paradigm to specify data analytics tasks.

Backend infrastructure as a highly-distributed execution environment for those tasks.

Performance. Largest Apache Spark cluster is 8000 nodes.

200 node Spark cluster sorted 100TB of data in 23 minutes.

Separation of concerns. Data analysts specifies analytics, doesn't worry about distributed nature of platform.



Programming Support for an **Integrated** Multi-Party Computation and MapReduce Infrastructure



To do big data analytics we need...

Data. But the interesting data is often private.

Performance. But MPC is slow!

An implementation of the analytics. But data analysts don't know MPC...

A concrete example

Goal: Want to establish the difference between male and female salaries across the big companies around Boston.

Good: Companies like this idea. Participating makes them look good.

Bad: Until they're asked to reveal their internal pay inequities to a “trusted” party.

If only they had our platform:

Each company can use **MapReduce** to find the salary differences in their own data.

→ Lots of computation

The companies can use **MPC** to find the collective difference without revealing their data.

→ Just one addition

The main components of our platform

Programming language to specify MapReduce and MPC operations.

Compiler to convert programs to tasks that are executable in existing MapReduce and MPC frameworks.

Backend platform running those MapReduce and MPC frameworks to act as an execution environment for a compiled program.

Let's explore our platform top-down

Programming language to specify MapReduce and MPC operations.

Compiler to convert programs to tasks that are executable in existing MapReduce and MPC frameworks.

Backend platform running those MapReduce and MPC frameworks to act as an execution environment for a compiled Scatter program.

Pay equity: declaring our **key-value** store

1: type gender = str

2: type salary = int

3: data := store(**gender**, *salary*)

Our MapReduce operations

1: type gender = str

2: type salary = int

3: data := **store**(gender, salary)

4: m := **reduce**(+, **filter**("m", data))

5: f := **reduce**(+, **filter**("f", data))

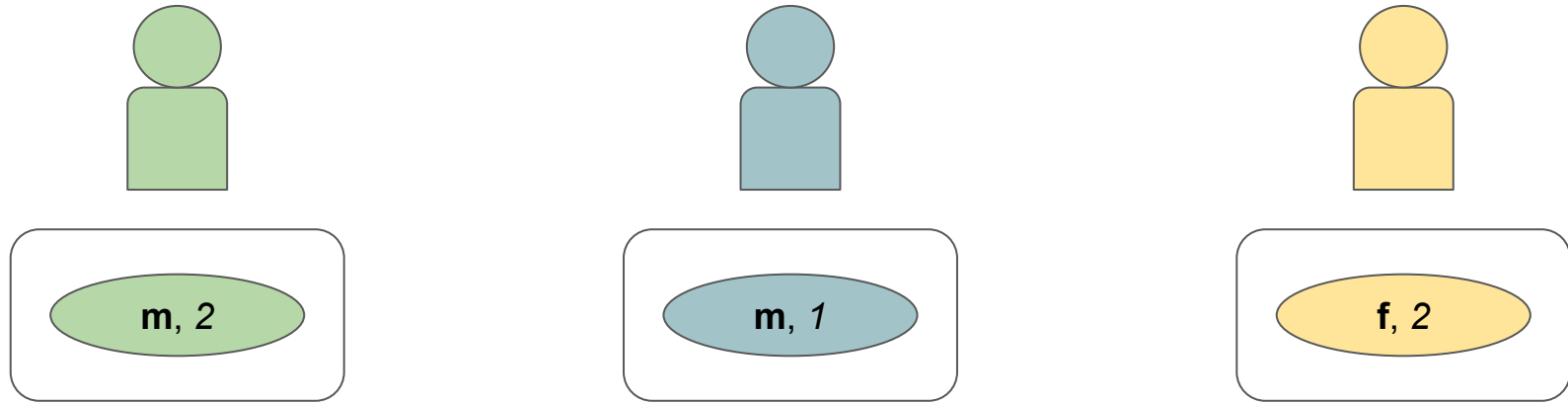
6: d := m - f

What about MPC?

Two main constructs:

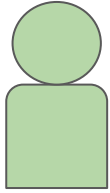
Scatter, and **Gather**.

Scatter: make secret and share



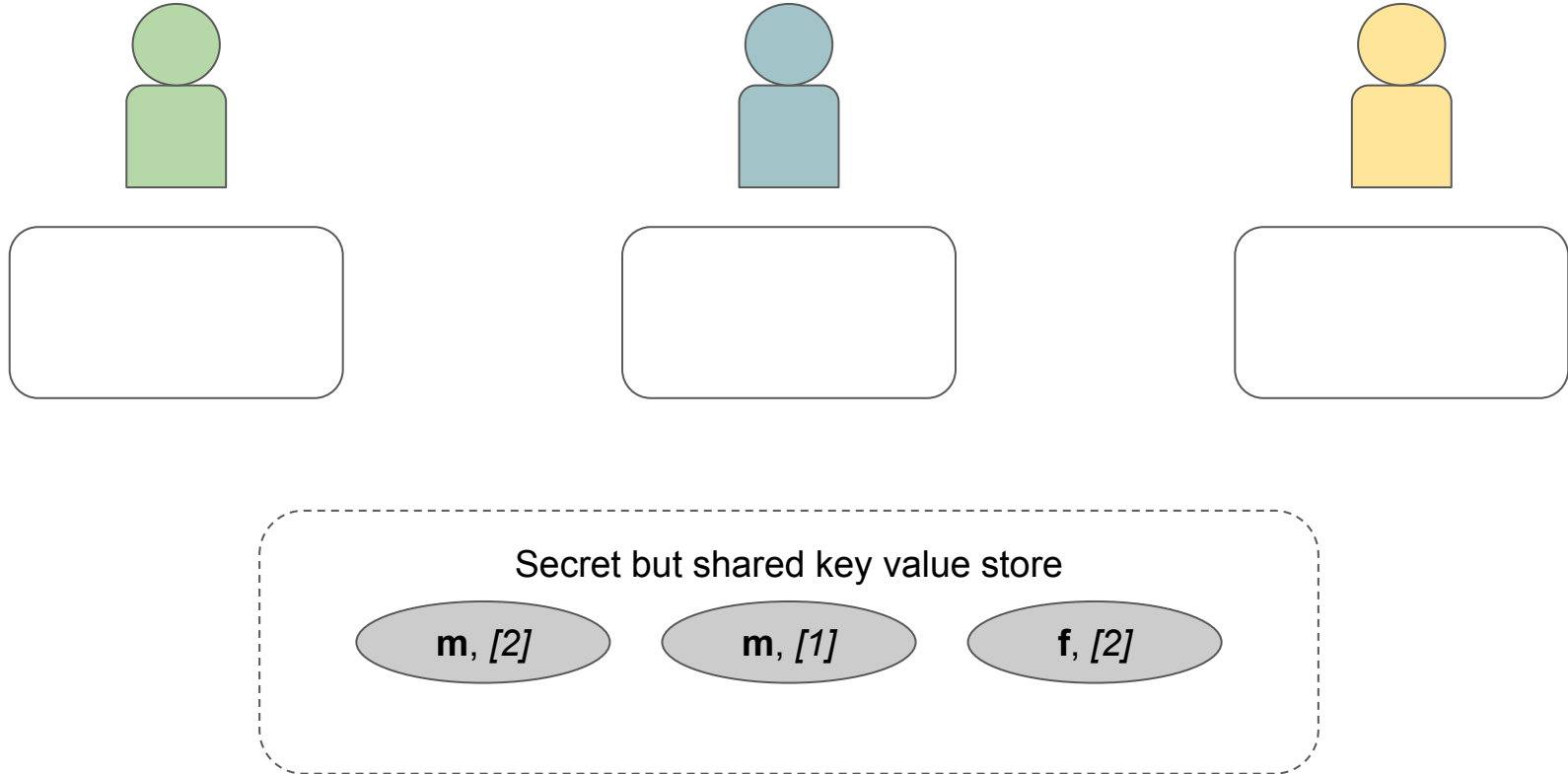
Secret but shared key value store

Scatter: make secret and share

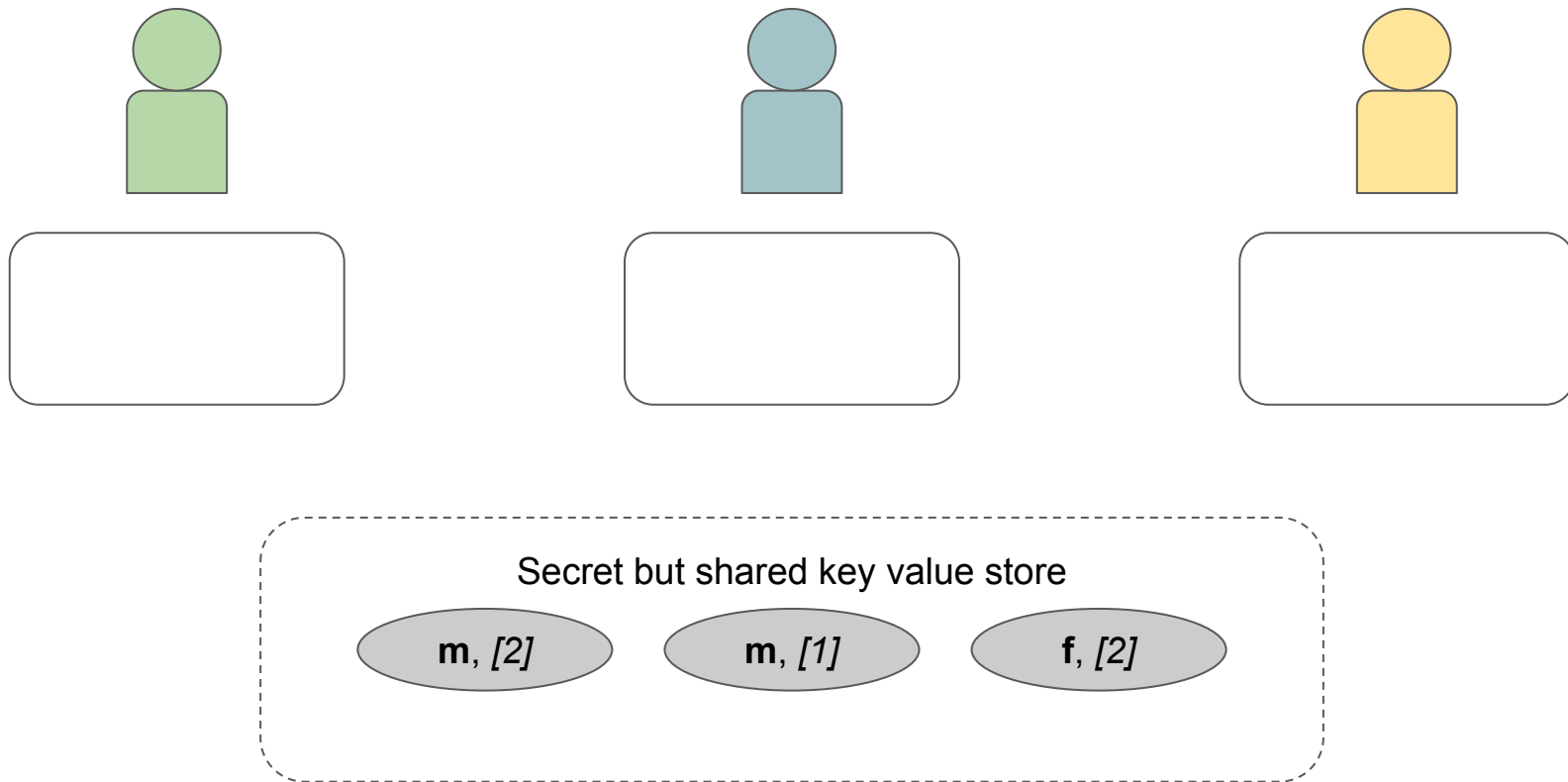


Secret but shared key value store

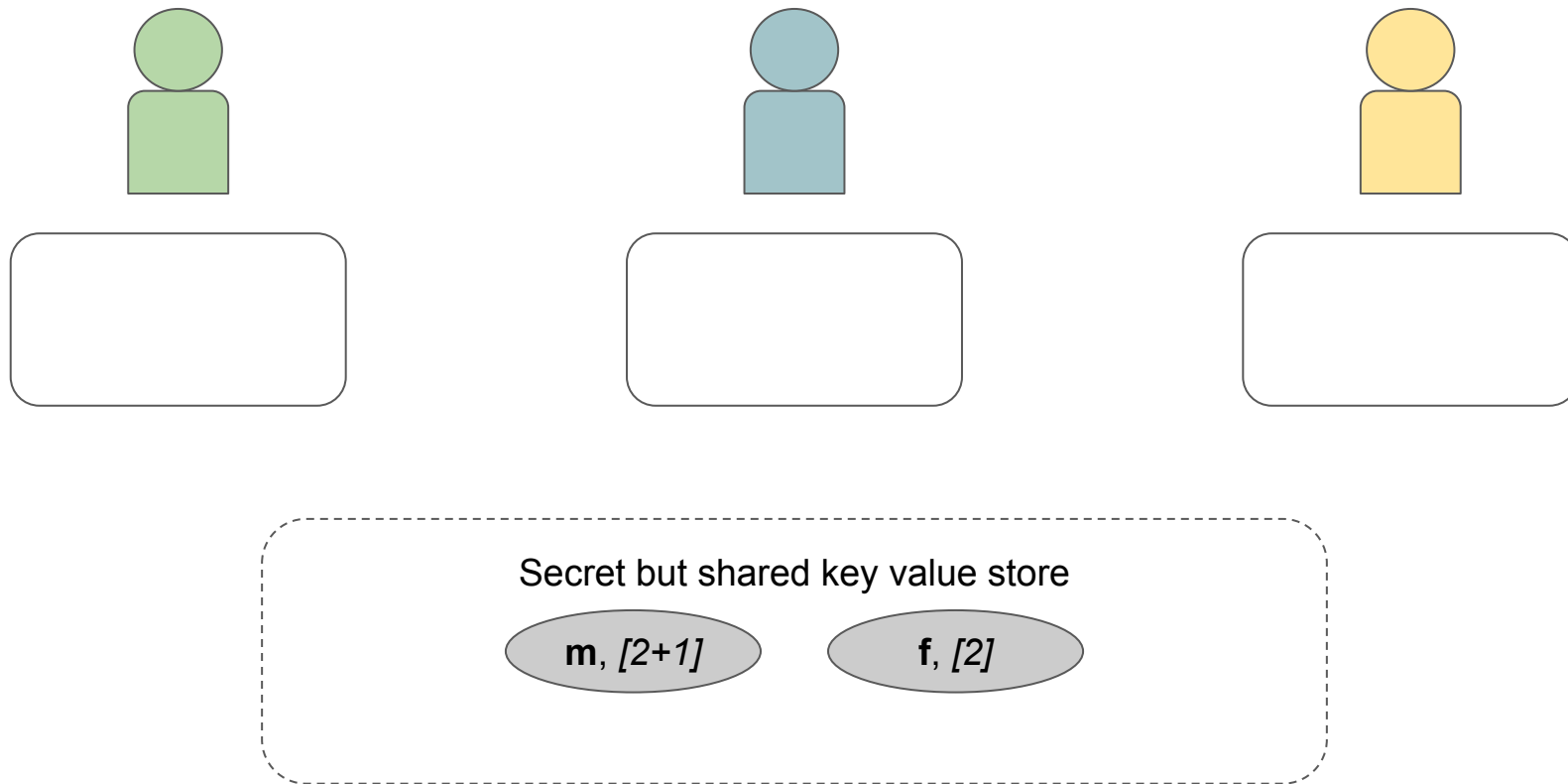
Scatter: make secret and share



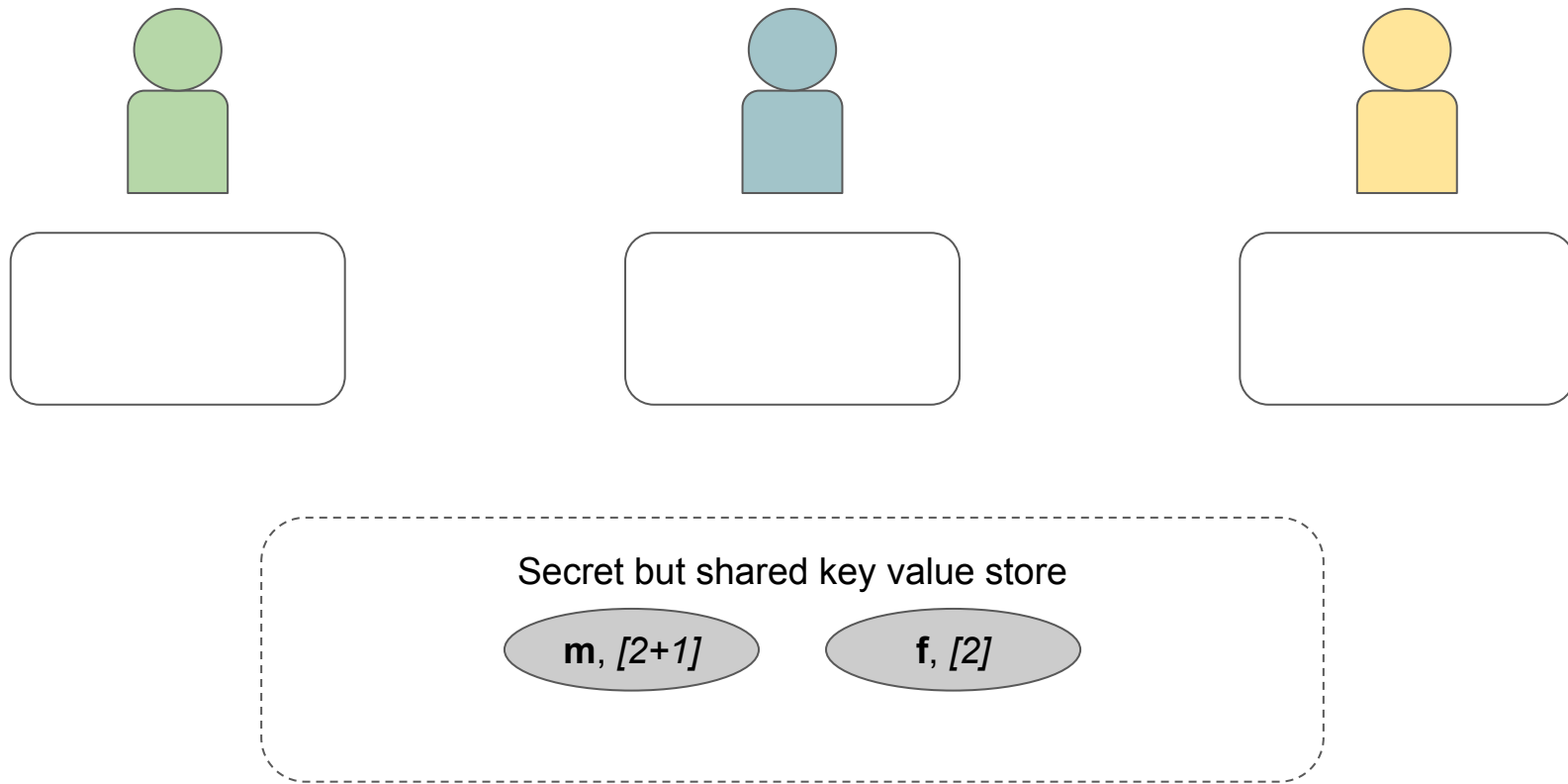
Let's say we want to perform a *reduce*(+)



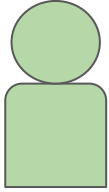
Let's say we want to perform a *reduce*(+)



Gather: collect and reveal

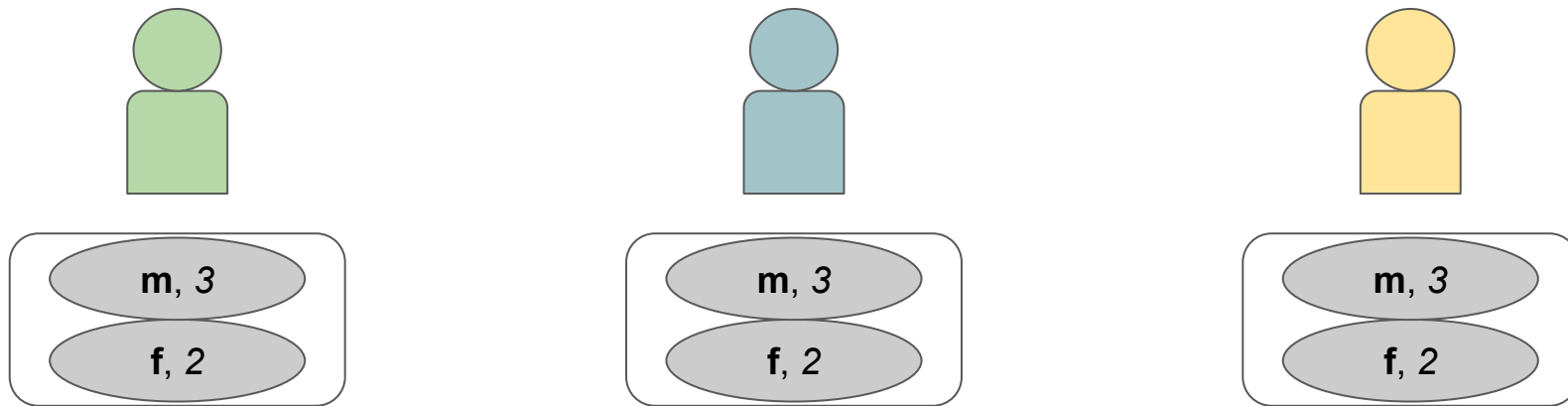


Gather: collect and reveal



Secret but shared key value store

Gather: collect and reveal



Secret but shared key value store

Complete pay equity program

```
1: type gender = str
2: type salary = int
3: data := store(gender, salary)

4: m := reduce(+, filter("m", data))
5: f := reduce(+, filter("f", data))
6: d := m - f

8: s := gather(reduce(+, scatter(d)))
```

Each company will execute this locally

1: type gender = str

2: type salary = int

3: data := store(gender, salary)

4: m := reduce(+, filter("m", data))

5: f := reduce(+, filter("f", data))

6: d := m - f

8: s := gather(reduce(+, scatter(d)))

The companies will need to perform an MPC

```
1: type gender = str
2: type salary = int
3: data := store(gender, salary)

4: m := reduce(+, filter("m", data))
5: f := reduce(+, filter("f", data))
6: d := m - f

8: s := gather(reduce(+, scatter(d)))
```

What do we do with a program?

Programming language to specify MapReduce and MPC operations.

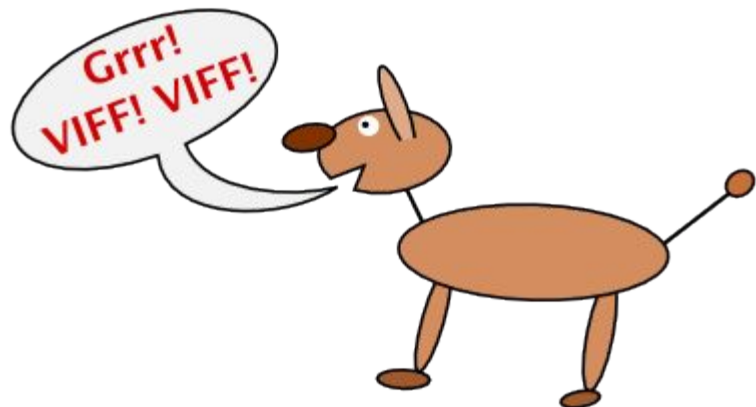
Compiler to convert Scatter programs to tasks that are executable in existing MapReduce and MPC frameworks.

Backend platform running those MapReduce and MPC frameworks to act as an execution environment for a compiled program.

Our current target frameworks



“a fast and general engine for large-scale data processing”



MPC framework that allows for Shamir secret sharing, arithmetic, and comparison over secret shares

4: $m := \text{reduce}(+, \text{filter}("m", \text{data}))$

5: $f := \text{reduce}(+, \text{filter}("f", \text{data}))$

6: $d := m - f$

```
m = data.filter(lambda x: x[0] == 'm')\
    .reduceByKey(lambda x, y: x + y)\
    .collect()
```

```
f = data.filter(lambda x: x[0] == 'f')\
    .reduceByKey(lambda x, y: x + y)\
    .collect()
```

```
d = ('d', m[0][1] - f[0][1])
```

8: $s := \text{gather}(\text{reduce}(+, \text{scatter}(d)))$

```
11  # Assume: 1. actor (shardid, k, v) == (id, k, kv_store)
12  mapped = map(lambda x: x[1], filtered)
13  private_kv_store = sorted(mapped, key=lambda x: x[0])
14
15  def protocol(rt):
16
17      def exchange_key_stores(rt):
18          def create_shared_key_stores(keys, player_mask):
19              keys_to_sharers = {}
20              [chr(k) for k in keys], player_mask
21              return keys_to_sharers
22
23          def key_store_sizes_ready(key_store_sizes):
24              return [int(k) for k in key_store_sizes]
25
26          def share_keys(key_store_sizes, rt, zp):
27              sorted_keys = []
28              player_mask = []
29              for player in rt.players:
30                  key_store_size = key_store_sizes[player - 1]
31                  for i in range(key_store_size):
32                      if rt.id == player:
33                          key = sorted(private_kv_store[i][0])
34                      else:
35                          key = None
36                      sorted_keys.append((rt.shamir_share(
37                          [player], zp, key, threshold=0))
38                          player_mask.append(player))
39
40              gathered_keys = gather_shares(
41                  [rt.open(k) for k in sorted_keys])
42
43              return gathered_keys.addCallback(create_shared_key_stores,
44                  player_mask)
45
46          shared_key_store_sizes = rt.shamir_share(
47              [players, zp, len(private_kv_store)], threshold=0)
48
49          opened_key_store_sizes = map(rt.open, shared_key_store_sizes)
50
51          key_store_sizes = gather_shares(
52              opened_key_store_sizes).addCallback(key_store_sizes_ready)
53
54          keys_to_sharers = key_store_sizes.addCallback(
55              share_keys, rt, zp)
56
57          return keys_to_sharers
58
59  def distribute_shares(rt, zp, private_kv_store, keys_to_sharers):
60      private_value_queue = collections.deque()
61      map(lambda x: x[1], private_kv_store)
62      shared_kv_store = []
63      for key, sharer in keys_to_sharers:
64          if sharer == rt.id:
65              value = rt.shamir_share(
66                  [sharer], zp, private_value_queue.popleft())
67              size:
68              value = rt.shamir_share([sharer], zp)
69              shared_kv_store.append((key, value))
70          return shared_kv_store
71
72  def open_shares(rt, kv_store, keys_to_owners, result_handler):
73      owner_queue = collections.deque(
74          map(lambda x: x[1], keys_to_owners))
75
76      opened_res = filter(lambda x: bool(x[1]), [(k, rt.open(v, owner_queue
77          for k, v in kv_store))]
78
79      expected_keys = sorted(map(lambda x: x[0], opened_res))
80      result_kv_store = []
81      for k, v in opened_res:
82          v.addCallback(
83              result_handler, k, result_kv_store, expected_keys)
84
85      rt.wait_for([map(lambda x: x[1], opened_res)])
86
87      def got_result(value, key, result_kv_store, expected_keys):
88          result_kv_store.append((key, int(value.sigmod))))
89          if expected_keys == sorted([x[0] for x in result_kv_store]):
90              output(result_kv_store, out_handle)
91
92      def continue_to_sharers():
93
```

What to do with executable code?

Programming language to (unified) specify MapReduce and MPC operations.

Compiler to convert programs to tasks that are executable in existing MapReduce and MPC frameworks.

Backend platform running those MapReduce and MPC frameworks to act as an execution environment for a program.

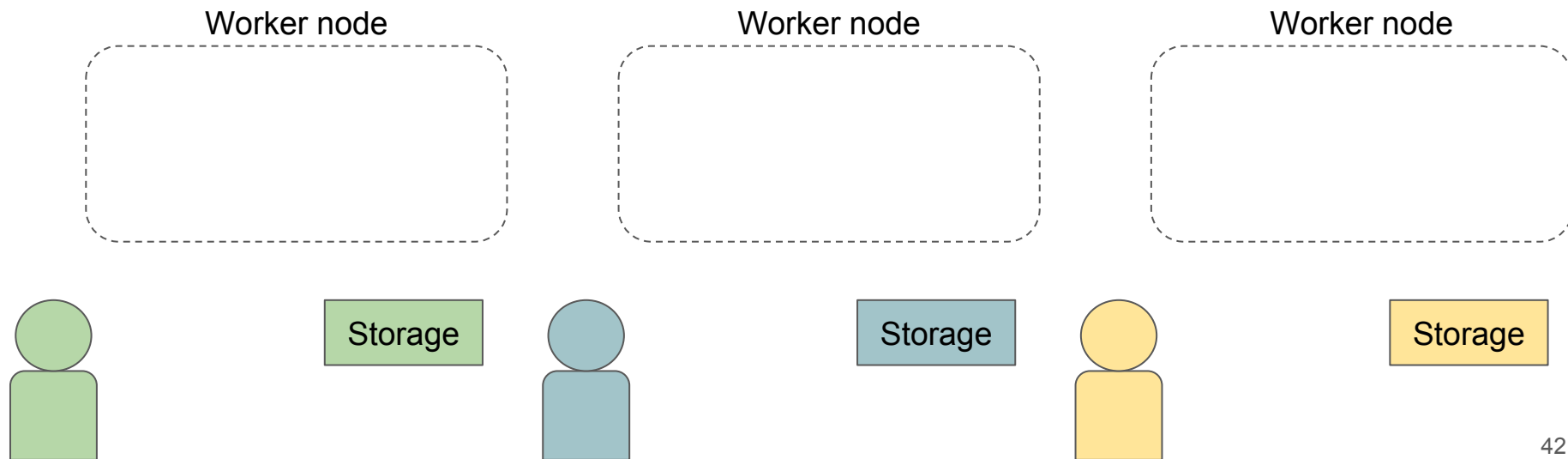
Let's build our backend.

Give each client the computational resources to:

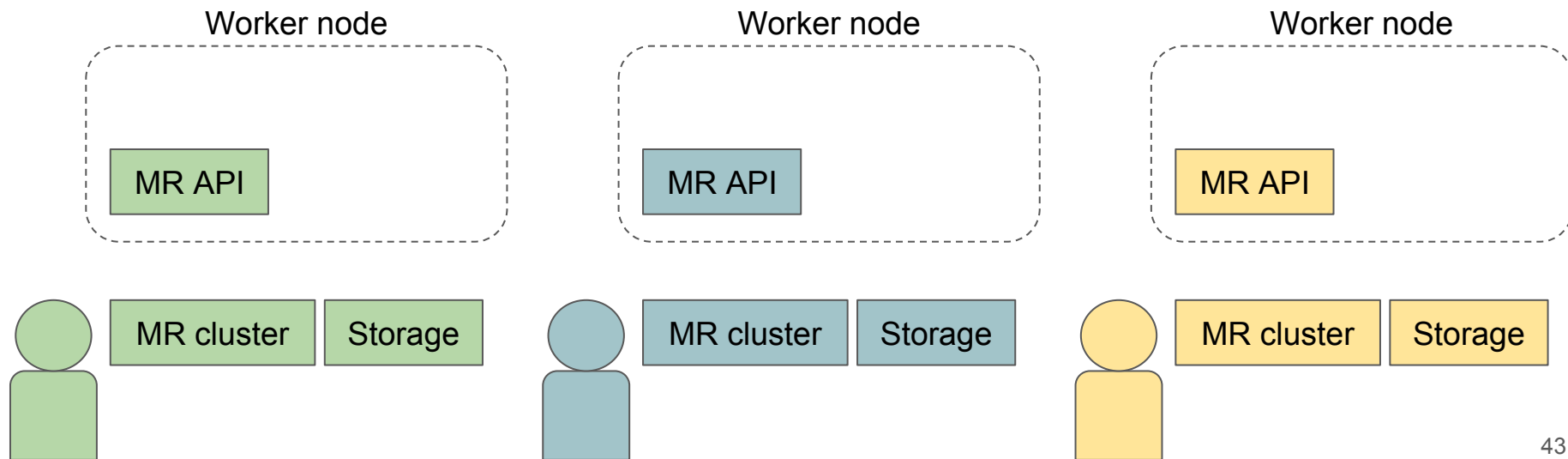
- run local MapReduce tasks on their data
- participate in MPC rounds to process data across companies
- coordinate those two actions

Let's build a **worker node**.

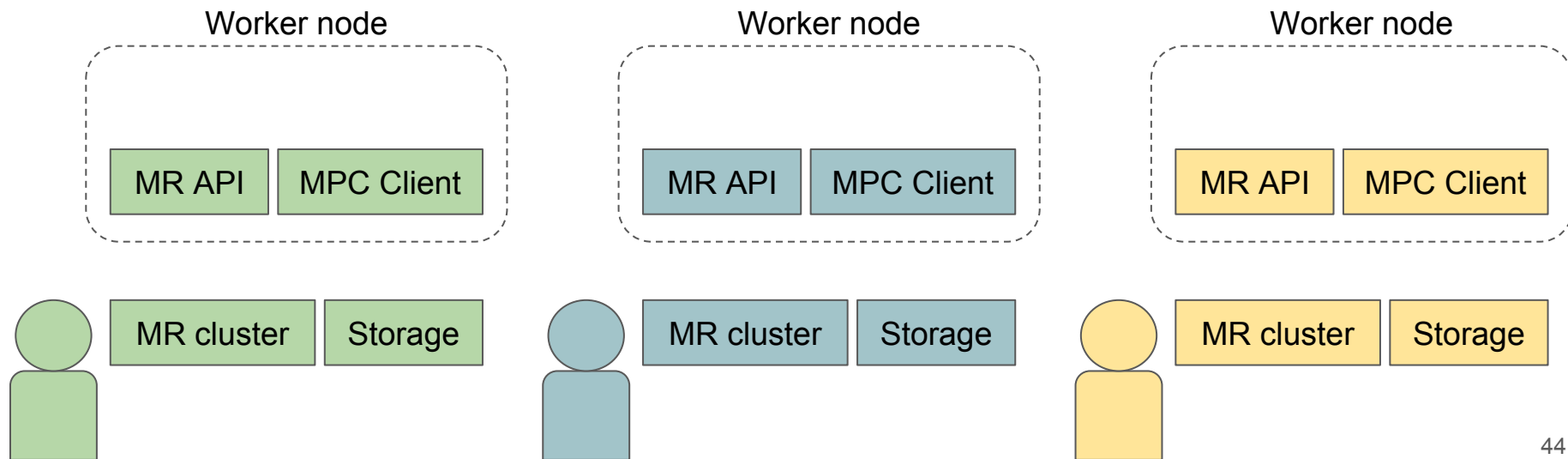
What does each company start with?



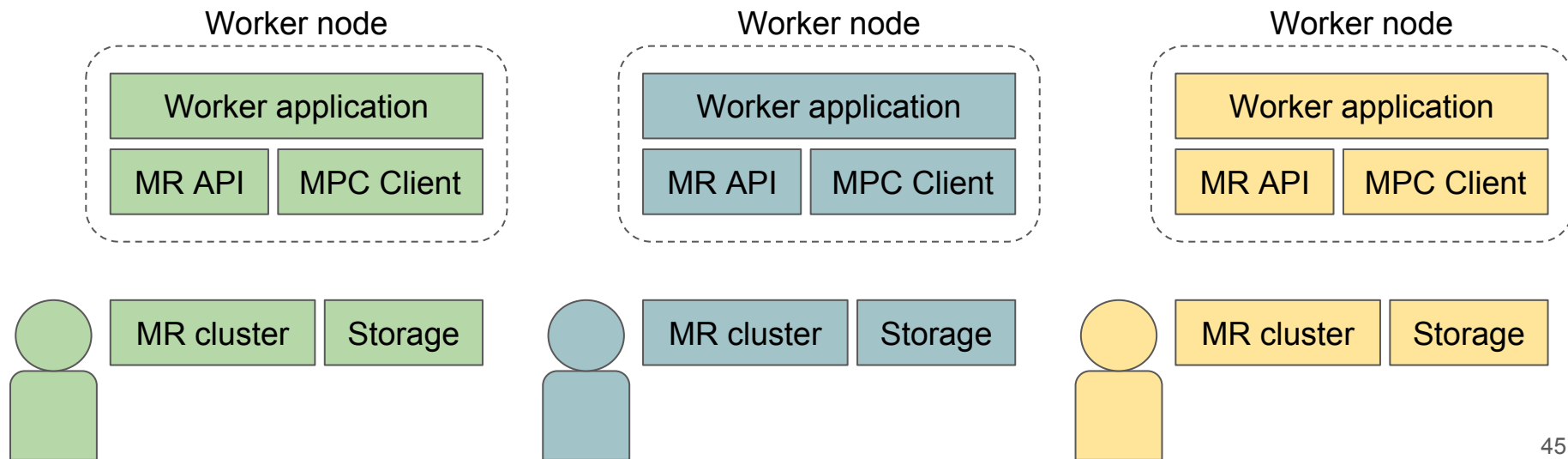
What do companies need to run MapReduce code?



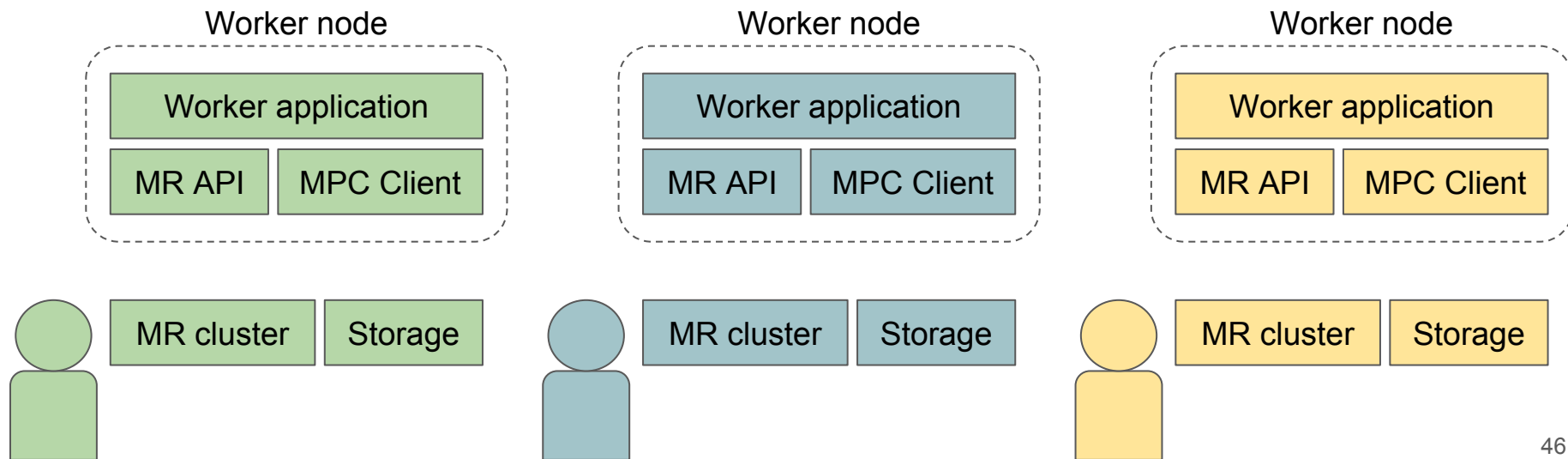
What do companies need to run MPC code?



What about coordinating program execution?



Bundle up the software, we have our **worker nodes**!



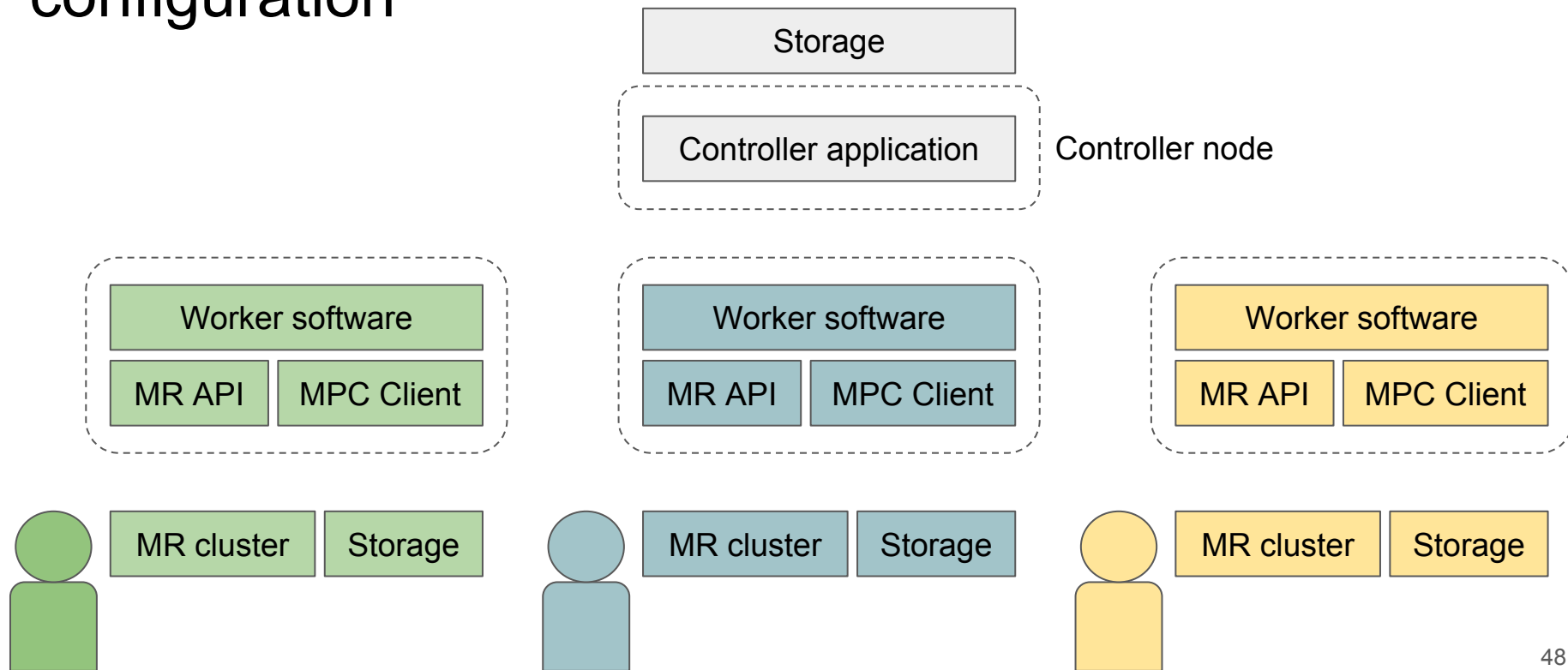
Almost done...

We have a distributed system.

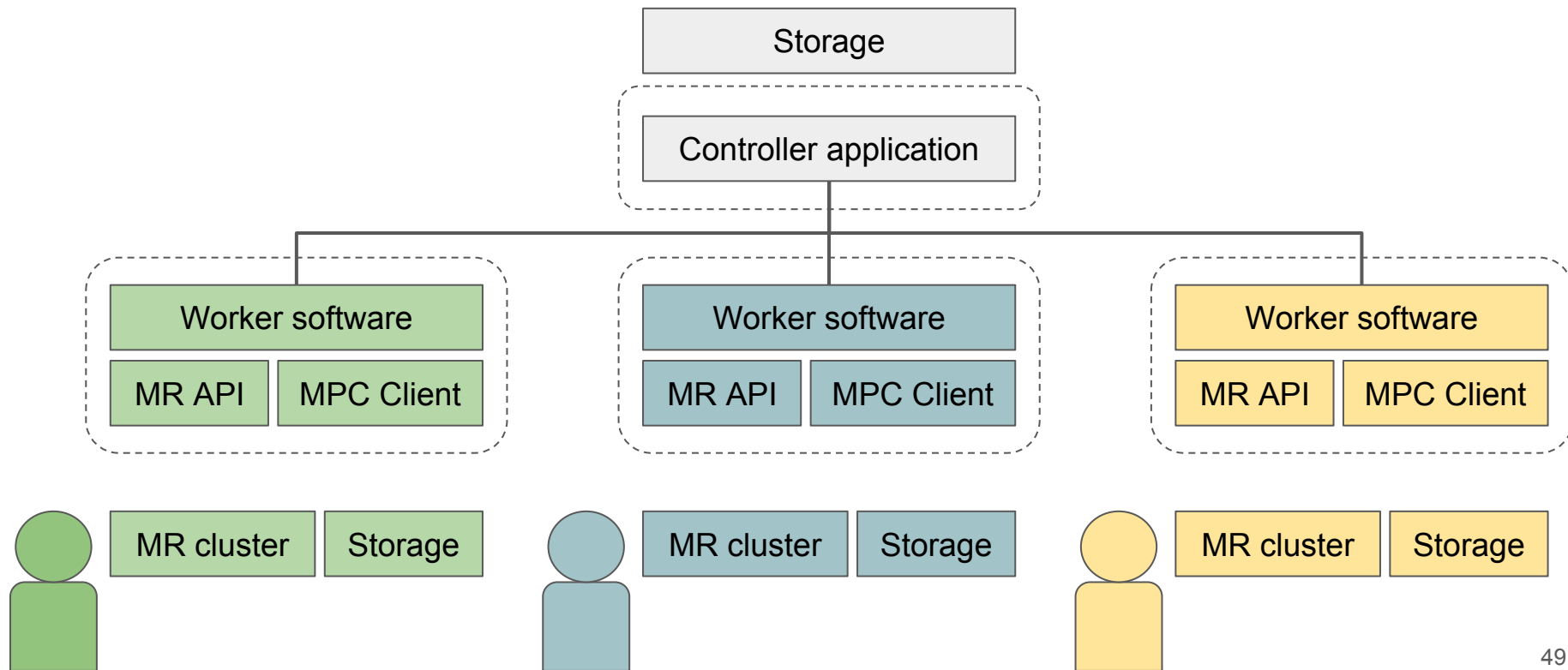
We need to coordinate task execution not only within worker nodes but also **across** worker nodes.

Let's build a **controller node**.

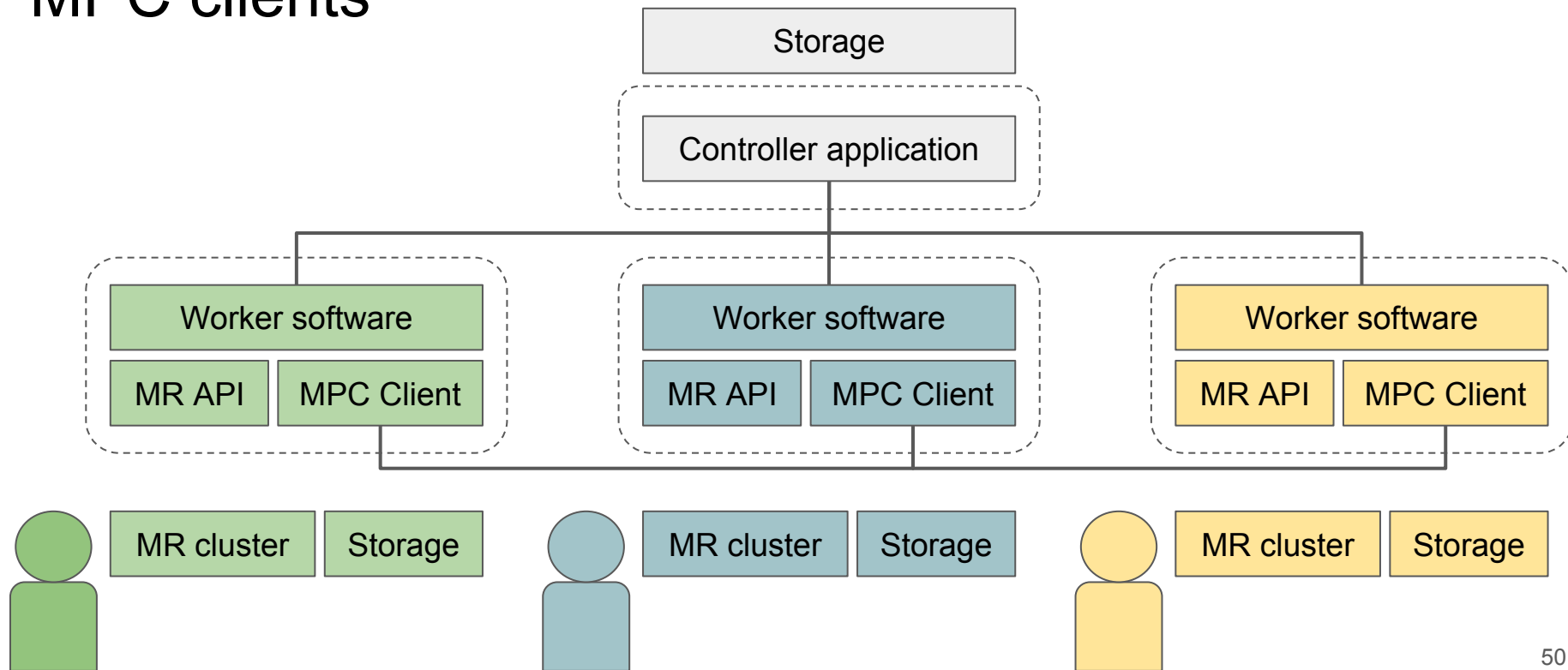
Controller: orchestrate task execution + worker configuration



Workers connect to Controller over HTTPS



Workers use controller to configure and connect MPC clients



And there we have it!

Programming language to specify MapReduce and MPC operations.

Compiler to convert programs to tasks that are executable in existing MapReduce and MPC frameworks.

Backend platform running those MapReduce and MPC frameworks to act as an execution environment for a compiled program.

Future work

Extending the support to more MPC and MapReduce backends

Separation of concerns

So many plans, so little time!

