

$$S = T_s / T_p$$

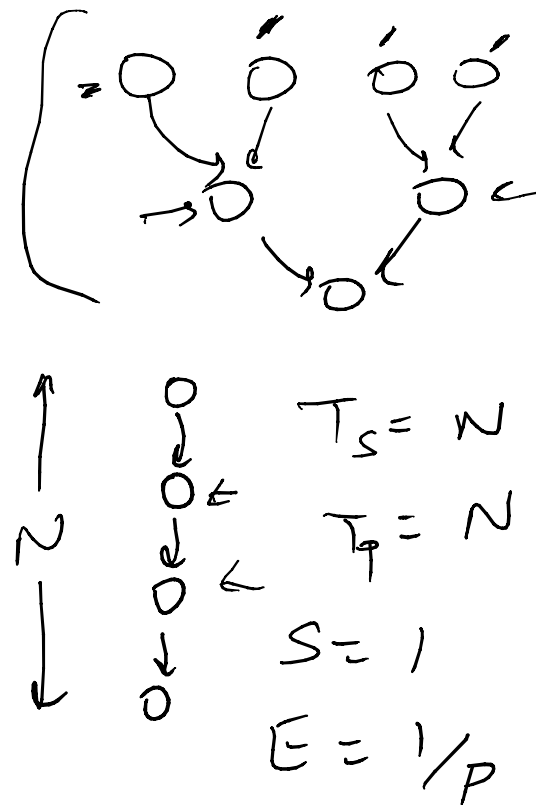
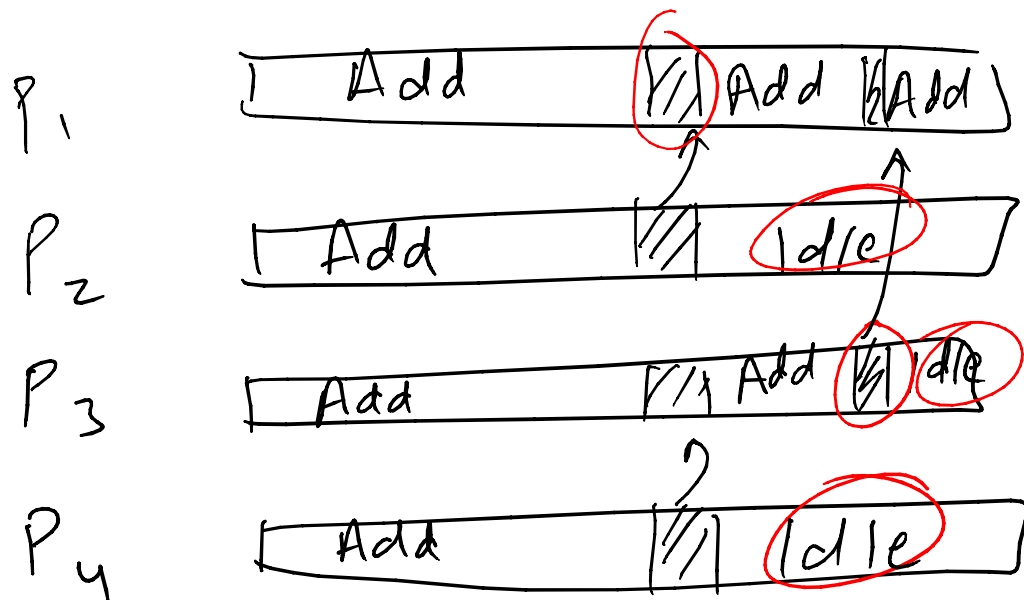
$$\frac{1.5}{3}$$

$$0.5$$

$$E = S/p = \frac{T_s}{p T_p}$$



$w =$



Implicit assumptions

- Uniform memory access
- No stragglers

Locality matters

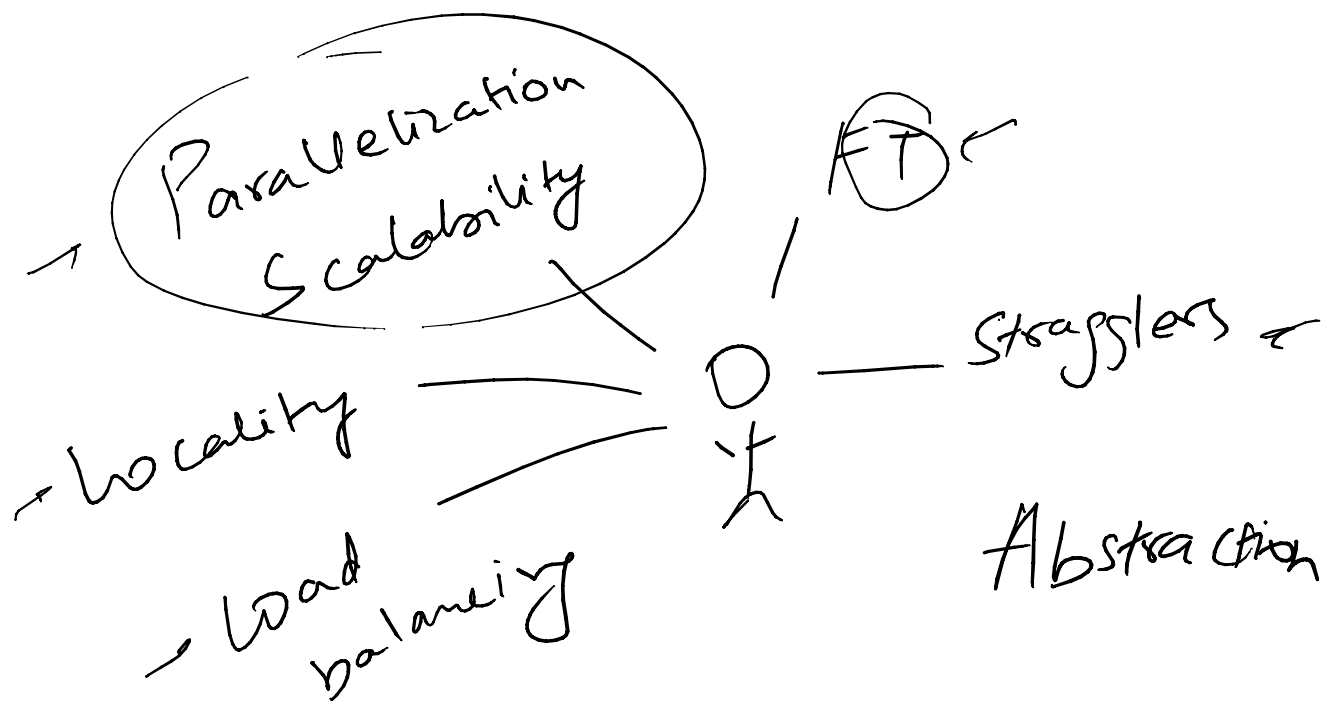
CPU \rightarrow cache $\sim 1\text{ns}$

DRAM $\sim 100\text{ns}$

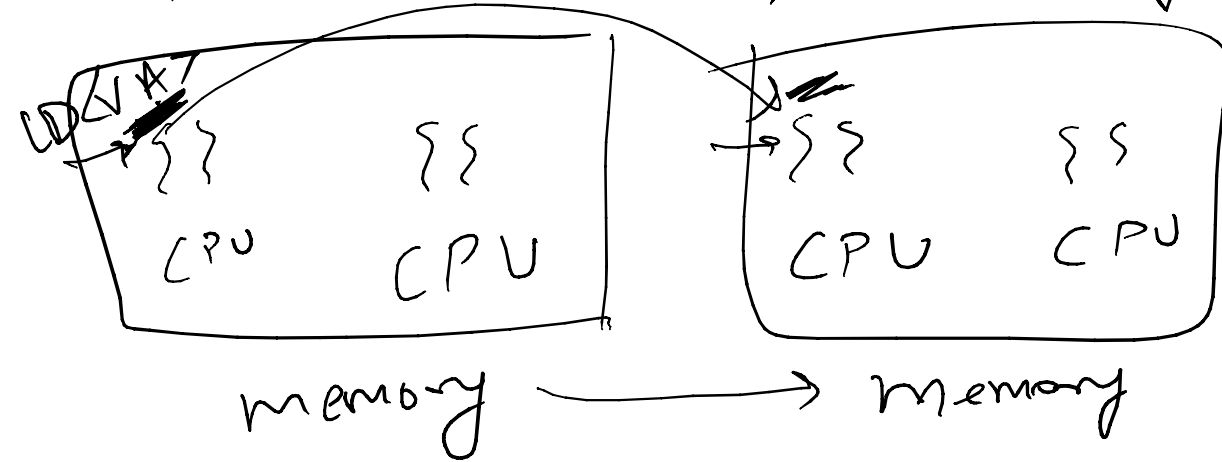
D/C $\sim 60\text{ps}$

Cross D/C $\sim 5-20\text{ns}$





Distributed Shared memory



Page-based DSM

TLB

VA	PA

→ Page table

VA	PA	Valid	Present

False sharing

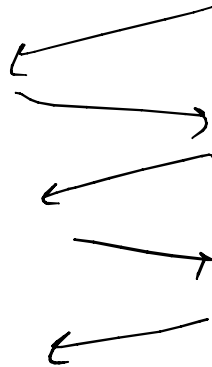
w1 for
 $\underline{x} \leftarrow x + 1$

w2 for

$y \leftarrow y + 1$

Page

x		
	y	



Who has the page?

Worker \longrightarrow

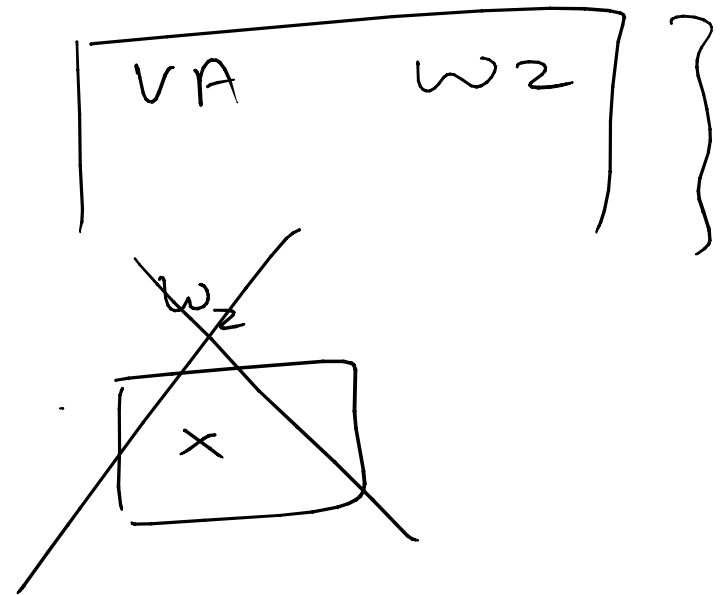
\swarrow bottleneck
(Central Server -

FT/Straggler

for

$x \leftarrow \underline{x} + 1$

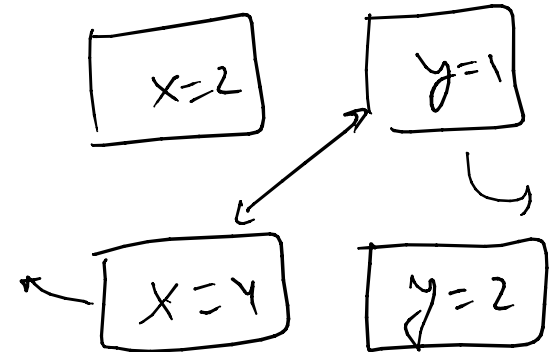
$\therefore \dots$



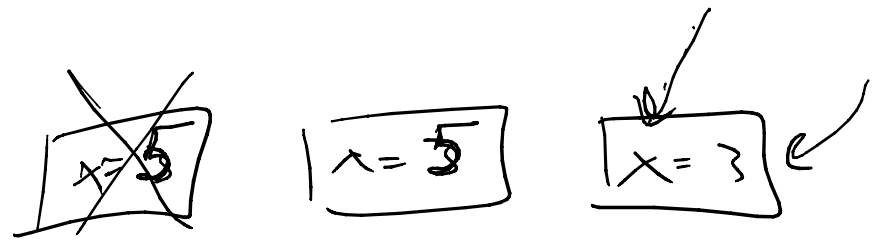
Fault Tolerance

- Checkpoint

- overhead
- regularly
- consistency
- Rollback



- Replication



Batch —

Streaming —

ML training —

Graph —

grep word
map { file, file contents }
for if Emit (line)
reduce:
append

Map Reduce

```
map(String key, String value):
```

```
// key: document name ←
```

```
// value: document contents ←
```

```
for each word w in value:
```

```
EmitIntermediate(w, "1"); ←
```

```
reduce(String key, Iterator values):
```

```
// key: a word ←
```

```
// values: a list of counts ←
```

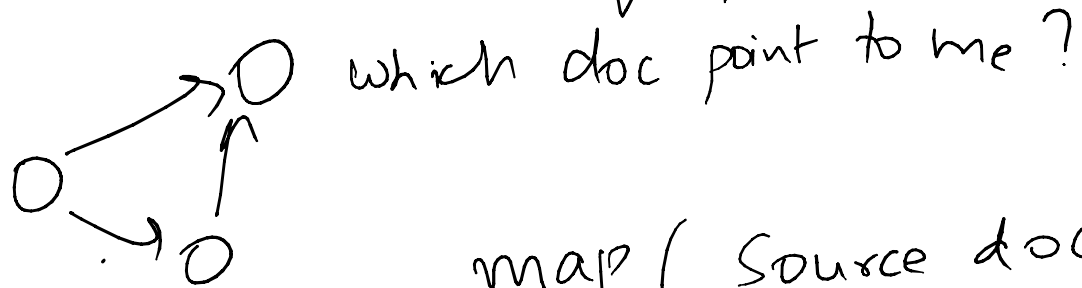
```
int result = 0;
```

```
for each v in values:
```

```
result += ParseInt(v); ←
```

```
Emit(AsString(result));
```

Reverse web link graph



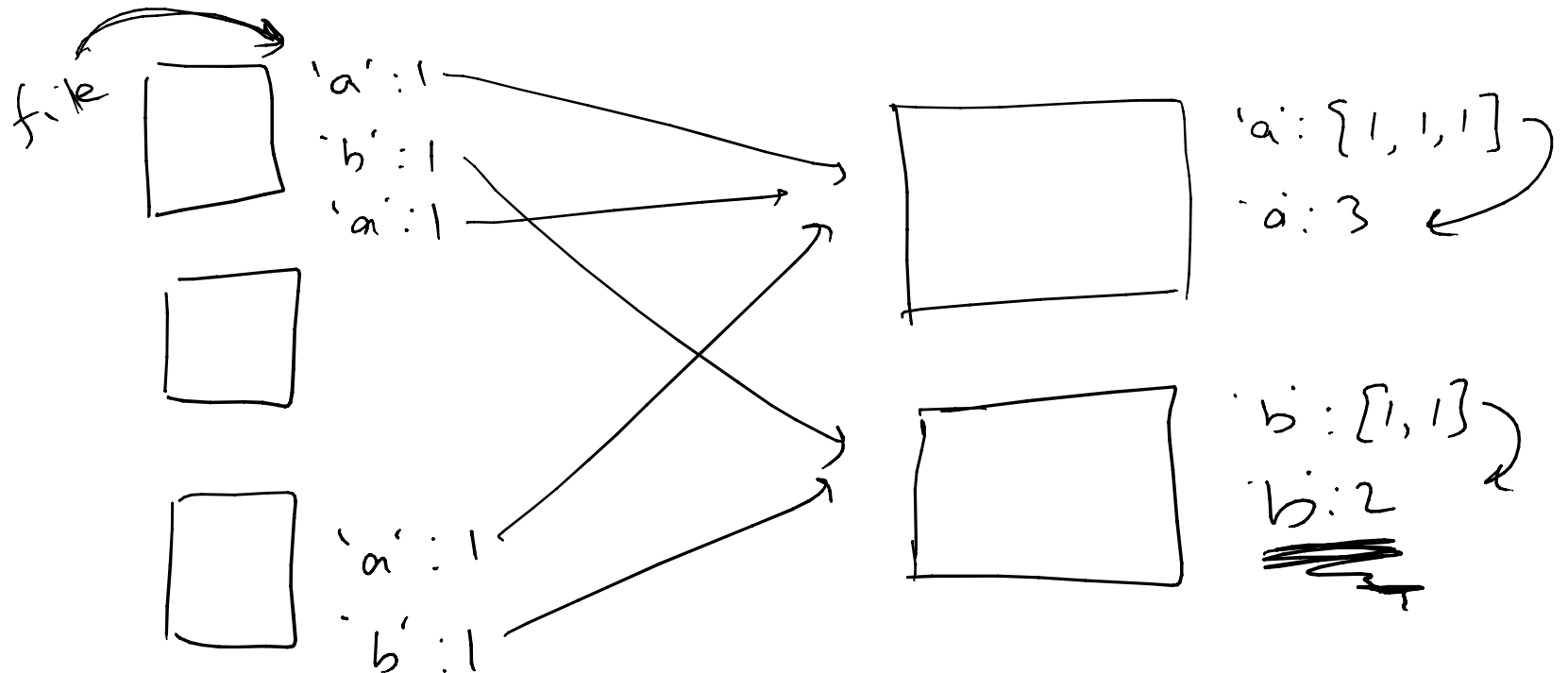
map (source doc)

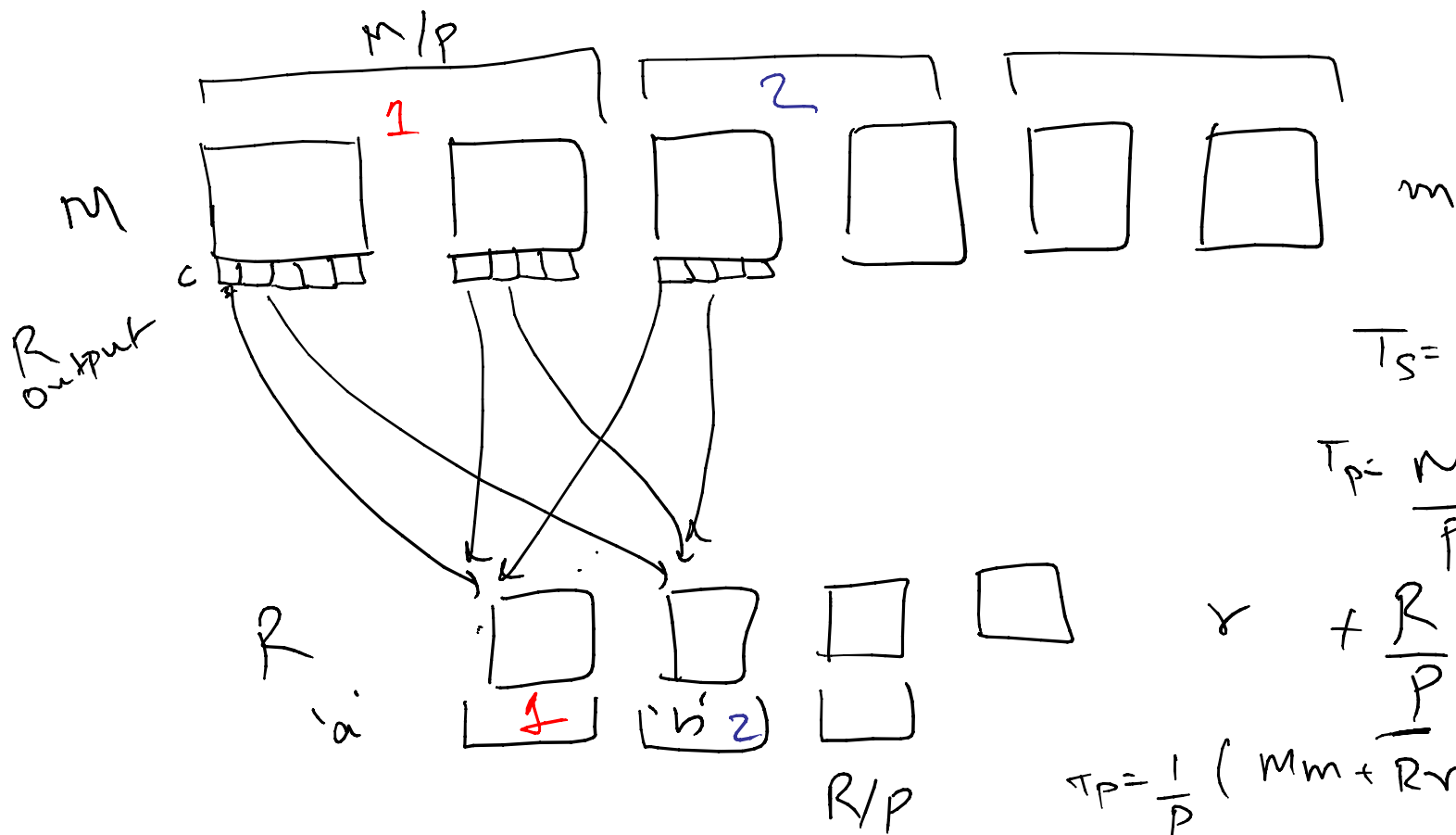
if link ^{parse} Emit (<target, source >)

Reduce (target,)

(target, list (sources))

Word count





$$T_s = Mm + Rr$$

$$T_p = \frac{M}{P}m + \frac{R}{P}r$$

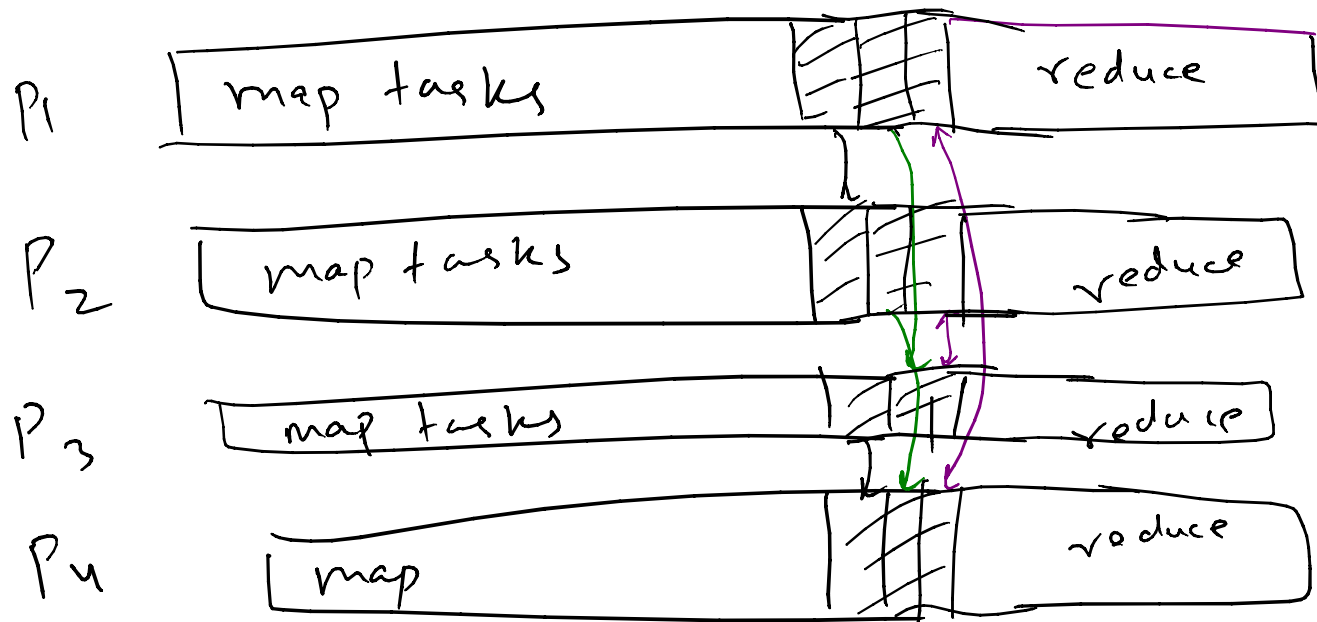
$$+ \frac{R}{P}(M - \frac{M}{P})$$

$$T_p = \frac{1}{P} (Mm + Rr + MRC \frac{(P-1)}{P})$$

$$T_p = \frac{1}{p} \left(M_m + R_r + MR \left(\frac{p-1}{p} \right) \right)$$

$$T_s = M_m + R_r$$

$$\underline{\underline{\epsilon}} = \frac{M_m + R_r}{M_m + R_r + MR \frac{(p-1)}{p}} = \frac{1}{1 + \left(\frac{MR \frac{(p-1)}{p}}{M_m + R_r} \right)}$$



$M, R \gg P$

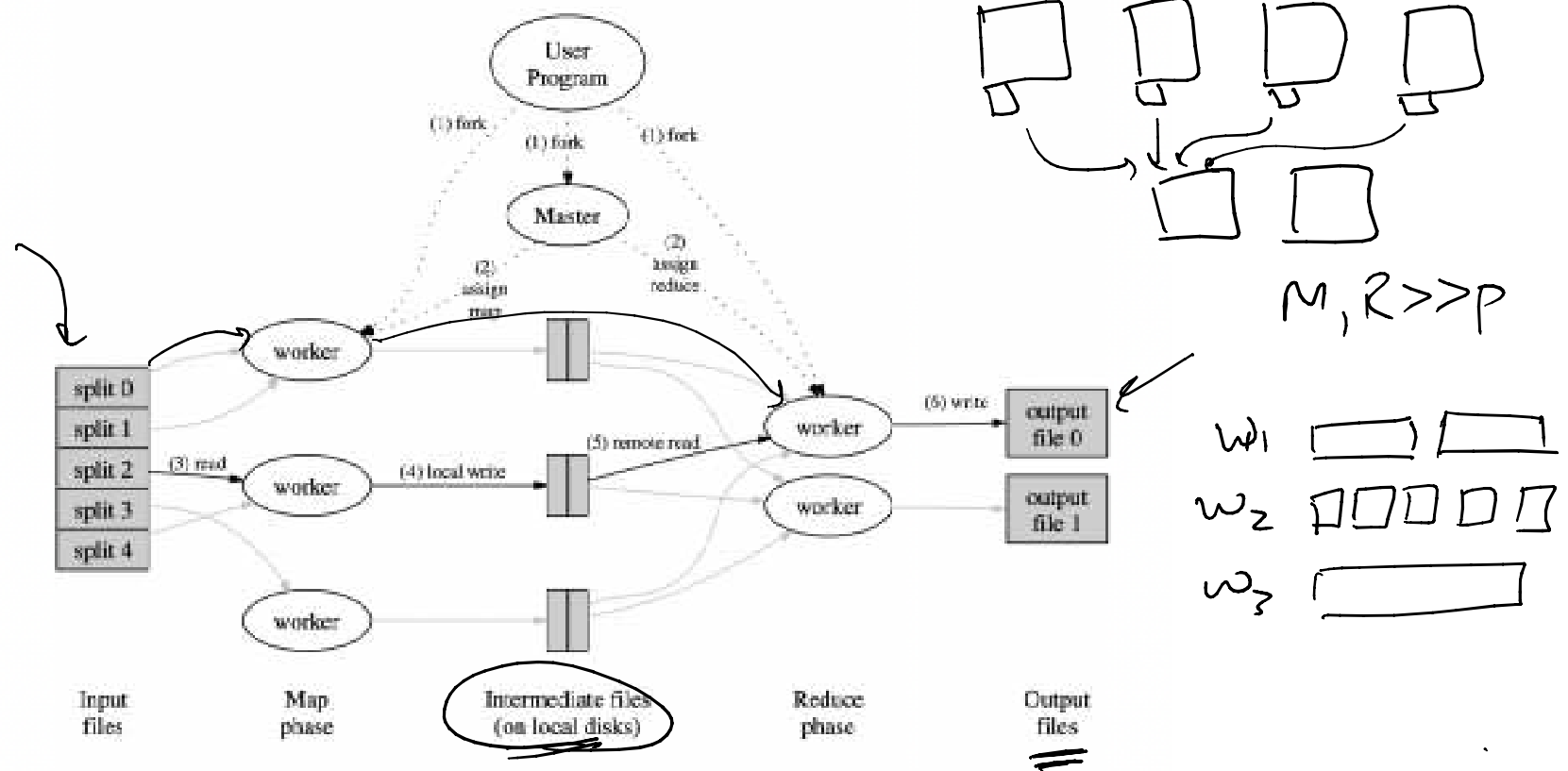


Figure 1: Execution overview

Master data structure

Tasks

m_1

m_2

Status

idle →

idle

in-process
↓
completed

worker

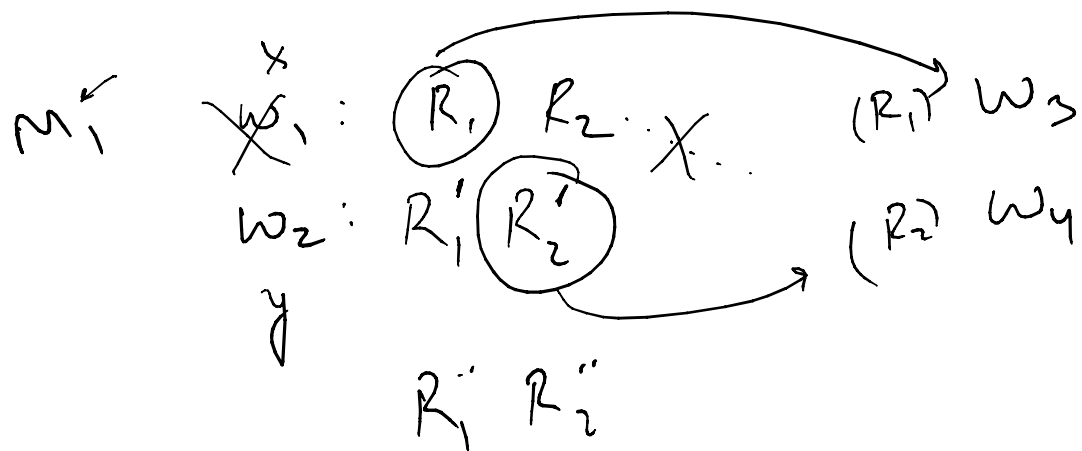
< R files, R sizes >

R_1

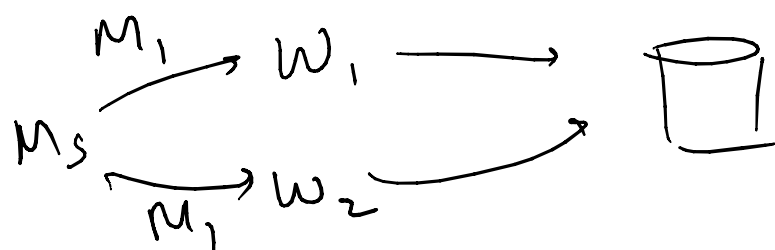
R_2

Fault tolerance -

- Re run task on a separate worker
- Heart Beats



$R_2 \neq R_2'$
deterministic tasks



Idempotence