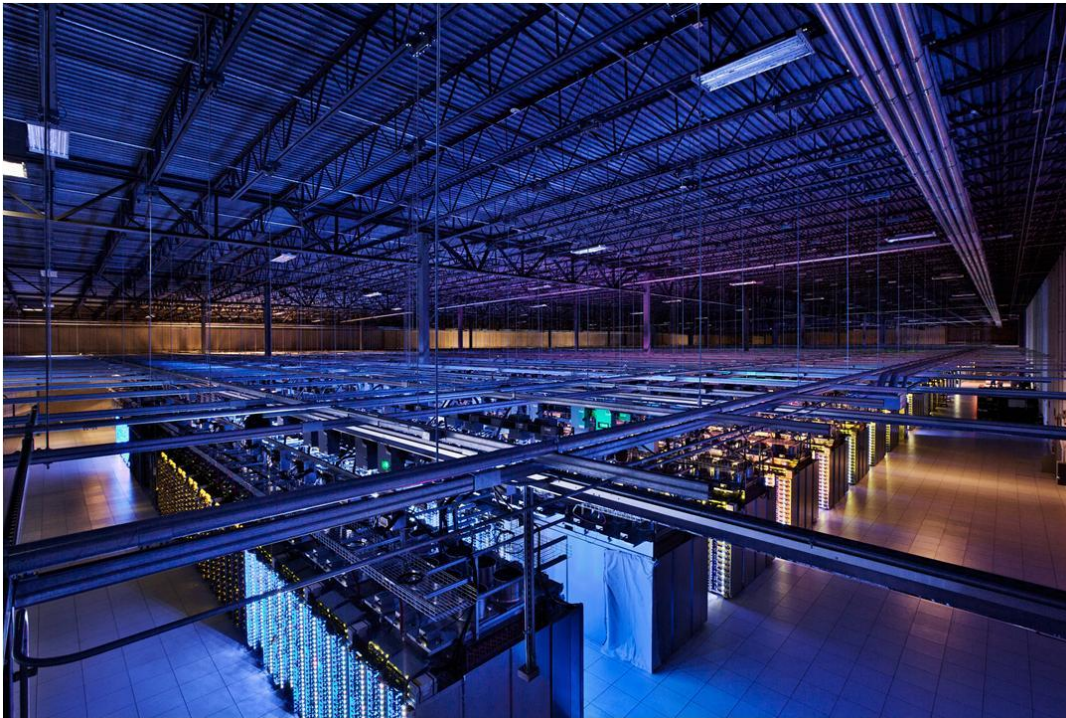


# Large-Scale Data Engineering

Designing and implementing algorithms  
for  
MapReduce

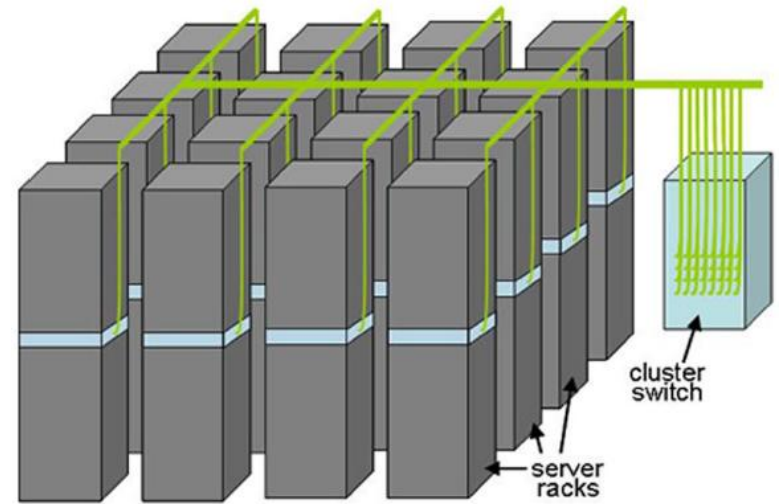
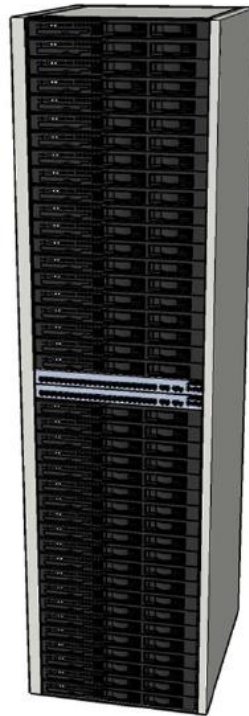
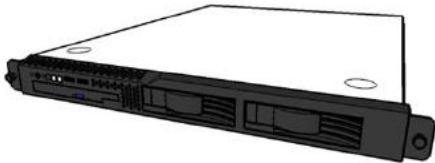


# PROGRAMMING FOR A DATA CENTRE

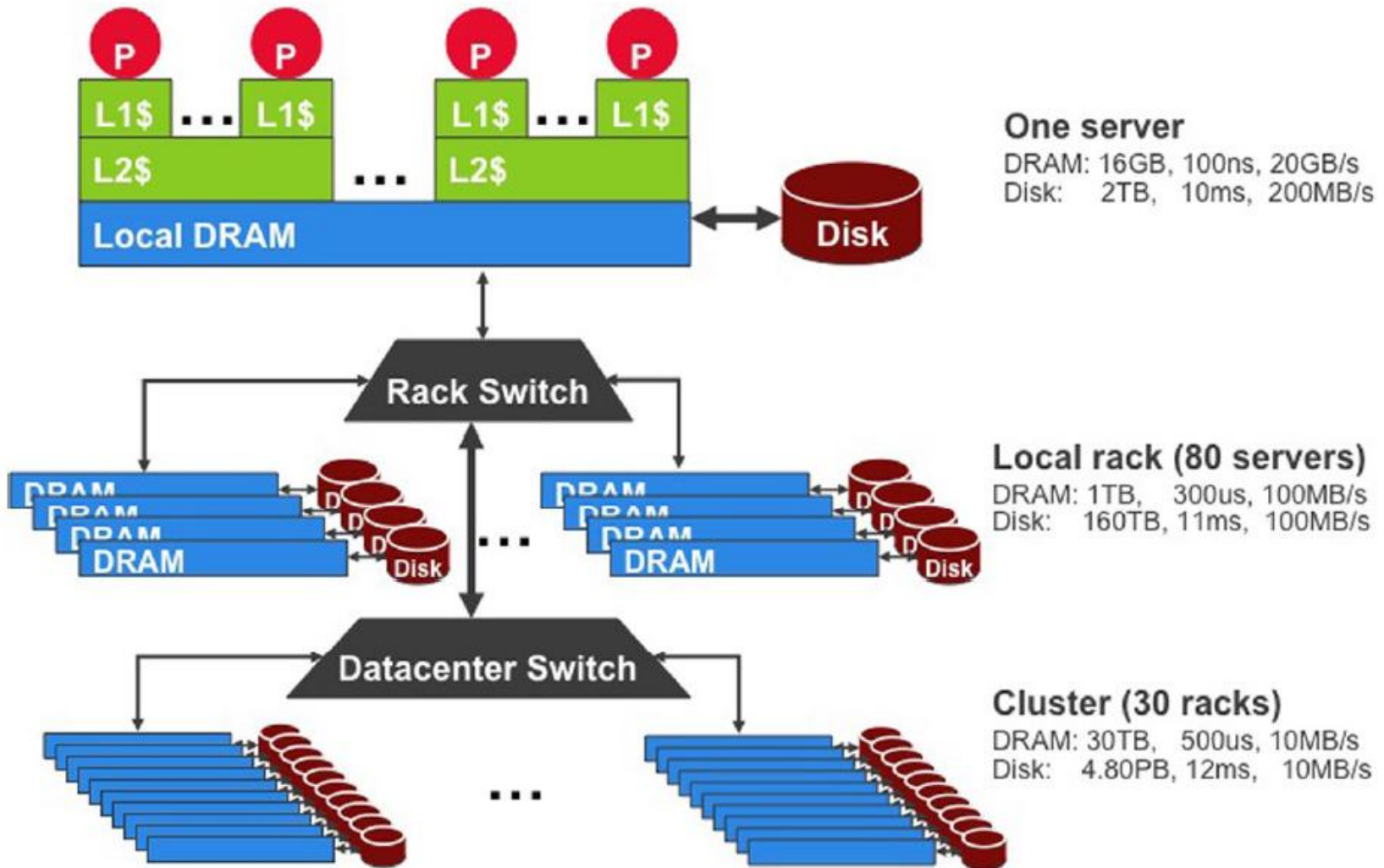
# Programming for a data centre

- Understanding the design of warehouse-sized computes
  - Different techniques for a different setting
  - Requires quite a bit of rethinking
- MapReduce algorithm design
  - How do you express everything in terms of `map()`, `reduce()`, `combine()`, and `partition()`?
  - Are there any design patterns we can leverage?

# Building Blocks



# Storage Hierarchy



# Scaling up vs. out

- No single machine is large enough
  - Smaller cluster of large SMP machines vs. larger cluster of commodity machines (e.g., 8 128-core machines vs. 128 8-core machines)
- Nodes need to talk to each other!
  - Intra-node latencies:  $\sim 100$  ns
  - Inter-node latencies:  $\sim 100$   $\mu$ s
- Let's model communication overhead

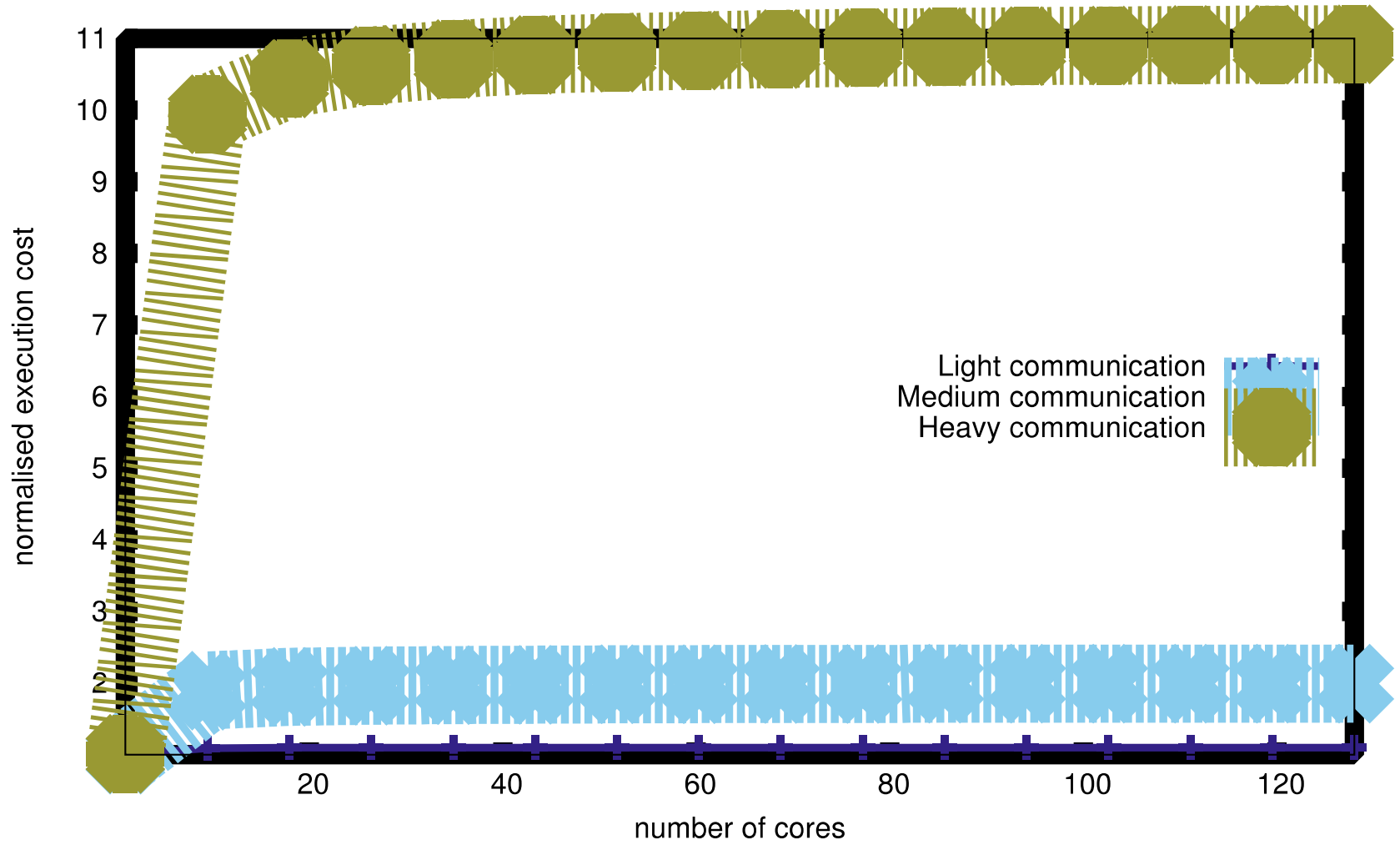
# Modelling communication overhead

- Simple execution cost model:
  - Total cost = cost of computation + cost to access global data
  - Fraction of local access inversely proportional to size of cluster
  - $n$  nodes (ignore cores for now)

$$1 \text{ ms} + f \times [100 \text{ ns} \times (1/n) + 100 \text{ } \mu\text{s} \times (1 - 1/n)]$$

- Light communication:  $f=1$
  - Medium communication:  $f=10$
  - Heavy communication:  $f=100$
- What is the cost of communication?

# Overhead of communication





# Seeks vs. scans

- Consider a 1TB database with 100 byte records
  - We want to update 1 percent of the records
- Scenario 1: random access
  - Each update takes ~30 ms (seek, read, write)
  - $10^8$  updates = ~35 days
- Scenario 2: rewrite all records
  - Assume 100MB/s throughput
  - Time = 5.6 hours(!)
- Lesson: avoid random seeks!

# Numbers everyone should know

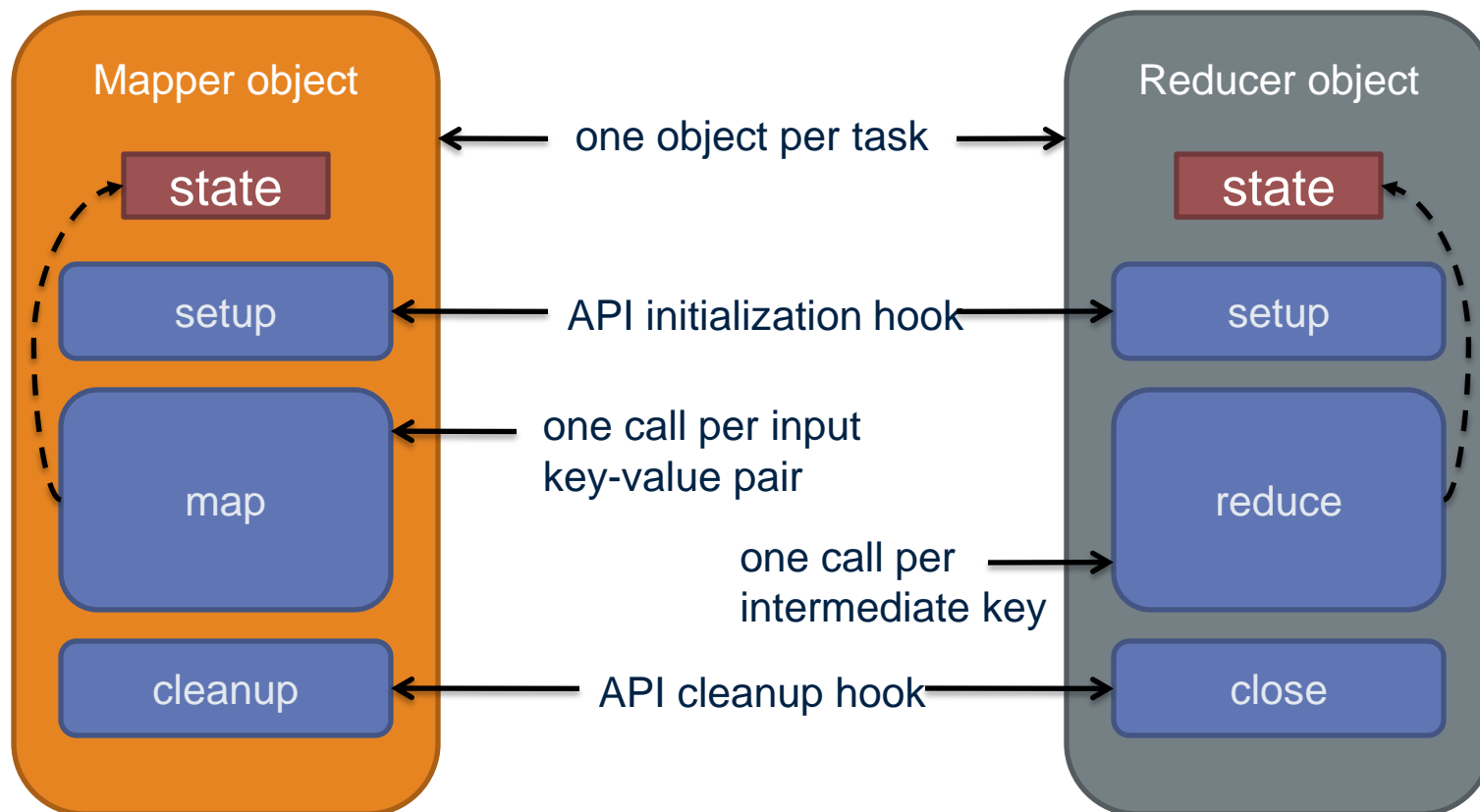
<b>L1 cache reference</b>	<b>0.5 ns</b>
Branch mispredict	5 ns
L2 cache reference	7 ns
Mutex lock/unlock	25 ns
Main memory reference	100 ns
Send 2K bytes over 1 Gbps network	20,000 ns
Read 1 MB sequentially from memory	250,000 ns
Round trip within same datacenter	500,000 ns
Disk seek	10,000,000 ns
Read 1 MB sequentially from disk	20,000,000 ns
Send packet CA → Netherlands → CA	150,000,000 ns

# DEVELOPING ALGORITHMS

# Optimising computation

- The cluster management software orchestrates the computation
- But we can still optimise the computation
  - Just as we can write better code and use better algorithms and data structures
  - At all times confined within the capabilities of the framework
- Cleverly-constructed data structures
  - Bring partial results together
- Sort order of intermediate keys
  - Control order in which reducers process keys
- Partitioner
  - Control which reducer processes which keys
- Preserving state in mappers and reducers
  - Capture dependencies across multiple keys and values

# Preserving State



# Importance of local aggregation

- Ideal scaling characteristics:
  - Twice the data, twice the running time
  - Twice the resources, half the running time
- Why can't we achieve this?
  - Synchronization requires communication
  - Communication kills performance
- Thus... avoid communication!
  - Reduce intermediate data via local aggregation
  - Combiners can help

# Word count: baseline

```
class Mapper
```

```
  method map(docid a, doc d)
```

```
    for all term t in d do
```

```
      emit(t, 1);
```

```
class Reducer
```

```
  method reduce(term t, counts [c1, c2, ...])
```

```
    sum = 0;
```

```
    for all counts c in [c1, c2, ...] do
```

```
      sum = sum + c;
```

```
    emit(t, sum);
```

# Word count: introducing combiners

```
class Mapper
  method map(docid a, doc d)
    H = associative_array(term → count;)
    for all term t in d do
      H[t]++;
    for all term t in H[t] do
      emit(t, H[t]);
```

Local aggregation reduces further computation



# Word count: introducing combiners

```
class Mapper
```

```
  method initialise()
```

```
    H = associative_array(term → count);
```

```
  method map(docid a, doc d)
```

```
    for all term t in d do
```

```
      H[t]++;
```

```
  method close()
```

```
    for all term t in H[t] do
```

```
      emit(t, H[t]);
```

Compute sums *across* documents!

# Design pattern for local aggregation

- In-mapper combining
  - Fold the functionality of the combiner into the mapper by preserving state across multiple map calls
- Advantages
  - Speed
  - Why is this faster than actual combiners?
- Disadvantages
  - Explicit memory management required
  - Potential for order-dependent bugs

# Combiner design

- Combiners and reducers share same method signature
  - Effectively they are map-side reducers
  - Sometimes, reducers can serve as combiners
  - Often, not...
- Remember: combiners are optional optimisations
  - Should not affect algorithm correctness
  - May be run 0, 1, or multiple times
- Example: find average of integers associated with the same key

# Computing the mean: version 1

```
class Mapper
```

```
  method map(string t, integer r)
```

```
    emit(t, r);
```

```
class Reducer
```

```
  method reduce(string, integers [r1, r2, ...])
```

```
    sum = 0;    count = 0;
```

```
    for all integers r in [r1, r2, ...] do
```

```
      sum = sum + r;    count++
```

```
     $r_{\text{avg}} = \text{sum} / \text{count};$ 
```

```
    emit(t,  $r_{\text{avg}}$ );
```

Can we use a reducer as the combiner?

# Computing the mean: version 2

```
class Mapper
```

```
  method map(string t, integer r)
```

```
    emit(t, r);
```

```
class Combiner
```

```
  method combine(string, integers [r1, r2, ...])
```

```
    sum = 0;    count = 0;
```

```
    for all integers r in [r1, r2, ...] do
```

```
      sum = sum + r;    count++;
```

```
      emit(t, pair(sum, count));
```

```
class Reducer
```

```
  method reduce(string, pairs [(s1, c1), (s2, c2), ...])
```

```
    sum = 0;    count = 0;
```

```
    for all pair(s, c) r in [(s1, c1), (s2, c2), ...] do
```

```
      sum = sum + s;    count = count + c;
```

```
    ravg = sum / count;
```

```
    emit(t, ravg);
```

Wrong!

# Computing the mean: version 3

```
class Mapper
```

```
  method map(string t, integer r)
```

```
    emit(t, pair(t, 1));
```

```
class Combiner
```

```
  method combine(string, pairs [(s1, c1), (s2, c2), ...])
```

```
    sum = 0;    count = 0;
```

```
    for all pair(s, c) in [(s1, c1), (s2, c2), ...] do
```

```
      sum = sum + s;    count = count + c;
```

```
    emit(t, pair(sum, count));
```

```
class Reducer
```

```
  method reduce(string, pairs [(s1, c1), (s2, c2), ...])
```

```
    sum = 0;    count = 0;
```

```
    for all pair(s, c) in [(s1, c1), (s2, c2), ...] do
```

```
      sum = sum + s;    count = count + c;
```

```
     $r_{\text{avg}} = \text{sum} / \text{count};$ 
```

```
    emit(t,  $r_{\text{avg}}$ );
```

Fixed!

# Computing the mean: version 4

```
class Mapper
```

```
  method initialise()
```

```
    S = associative_array(string → integer);
```

```
    C = associative_array(string → integer);
```

```
  method map(string t, integer r)
```

```
    S[t] = S[t] + r;    C[t]++;
```

```
  method close()
```

```
    for all t in keys(S) do
```

```
      emit(t, pair(S[t], C[t]));
```

Simpler, cleaner, with no need for combiner

# Algorithm design: term co-occurrence

- Term co-occurrence matrix for a text collection
  - $M = N \times N$  matrix ( $N =$  vocabulary size)
  - $M_{ij}$ : number of times  $i$  and  $j$  co-occur in some context (for concreteness, let's say context = sentence)
- Why?
  - Distributional profiles as a way of measuring semantic distance
  - Semantic distance useful for many language processing tasks



# Using MapReduce for large counting problems

- Term co-occurrence matrix for a text collection is a specific instance of a large counting problem
  - A large event space (number of terms)
  - A large number of observations (the collection itself)
  - Goal: keep track of interesting statistics about the events
- Basic approach
  - Mappers generate partial counts
  - Reducers aggregate partial counts

How do we aggregate partial counts efficiently?

# First try: pairs

- Each mapper takes a sentence:
  - Generate all co-occurring term pairs
  - For all pairs, emit  $(a, b) \rightarrow \text{count}$
- Reducers sum up counts associated with these pairs
- Use combiners!

# Pairs: pseudo-code

```
class Mapper
```

```
  method map(docid a, doc d)
```

```
    for all w in d do
```

```
      for all u in neighbours(w) do
```

```
        emit(pair(w, u), 1);
```

```
class Reducer
```

```
  method reduce(pair p, counts [c1, c2, ...])
```

```
    sum = 0;
```

```
    for all c in [c1, c2, ...] do
```

```
      sum = sum + c;
```

```
    emit(p, sum);
```

# Analysing pairs

- Advantages
  - Easy to implement, easy to understand
- Disadvantages
  - Lots of pairs to sort and shuffle around (upper bound?)
  - Not many opportunities for combiners to work

# Another try: stripes

- Idea: group together pairs into an associative array

$(a, b) \rightarrow 1$

$(a, c) \rightarrow 2$

$(a, d) \rightarrow 5$

$(a, e) \rightarrow 3$

$(a, f) \rightarrow 2$

$a \rightarrow \{ b: 1, c: 2, d: 5, e: 3, f: 2 \}$

- Each mapper takes a sentence:
  - Generate all co-occurring term pairs
  - For each term, emit  $a \rightarrow \{ b: \text{count}_b, c: \text{count}_c, d: \text{count}_d \dots \}$
- Reducers perform element-wise sum of associative arrays

$a \rightarrow \{ b: 1, \quad d: 5, e: 3 \}$

$a \rightarrow \{ b: 1, c: 2, \quad d: 2, \quad f: 2 \}$

$a \rightarrow \{ b: 2, c: 2, \quad d: 7, e: 3, \quad f: 2 \}$

Cleverly-constructed data structure brings together partial results

# Stripes: pseudo-code

```
class Mapper
```

```
  method map(docid a, doc d)
```

```
    for all w in d do
```

```
      H = associative_array(string → integer);
```

```
      for all u in neighbours(w) do
```

```
        H[u]++;
```

```
      emit(w, H);
```

```
class Reducer
```

```
  method reduce(term w, stripes [H1, H2, ...])
```

```
    Hf = associative_array(string → integer);
```

```
    for all H in [H1, H2, ...] do
```

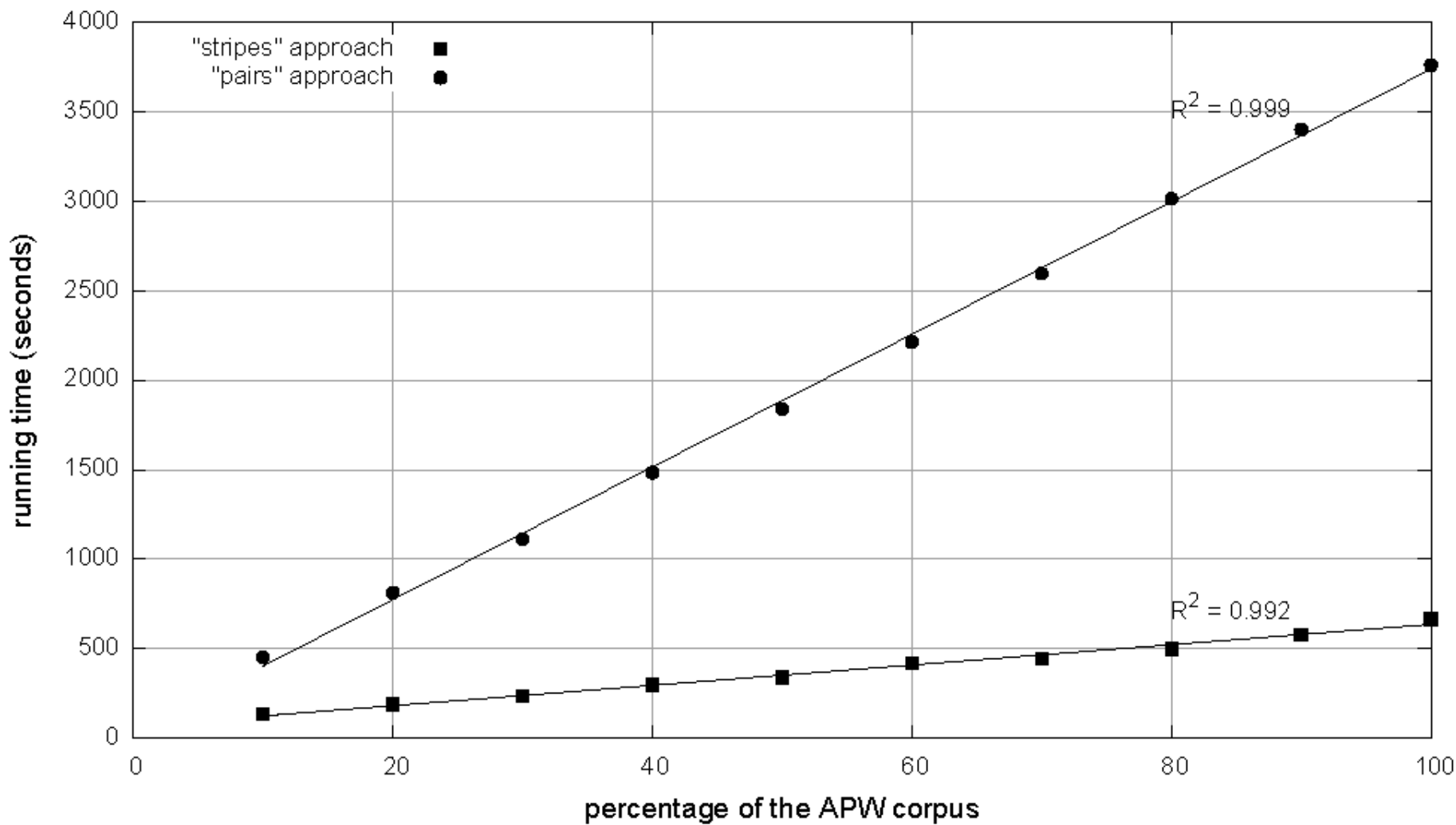
```
      sum(Hf, H);    // sum same-keyed entries
```

```
    emit(w, Hf);
```

# Stripes analysis

- Advantages
  - Far less sorting and shuffling of key-value pairs
  - Can make better use of combiners
- Disadvantages
  - More difficult to implement
  - Underlying object more heavyweight
  - Fundamental limitation in terms of size of event space

## Comparison of "pairs" vs. "stripes" for computing word co-occurrence matrices

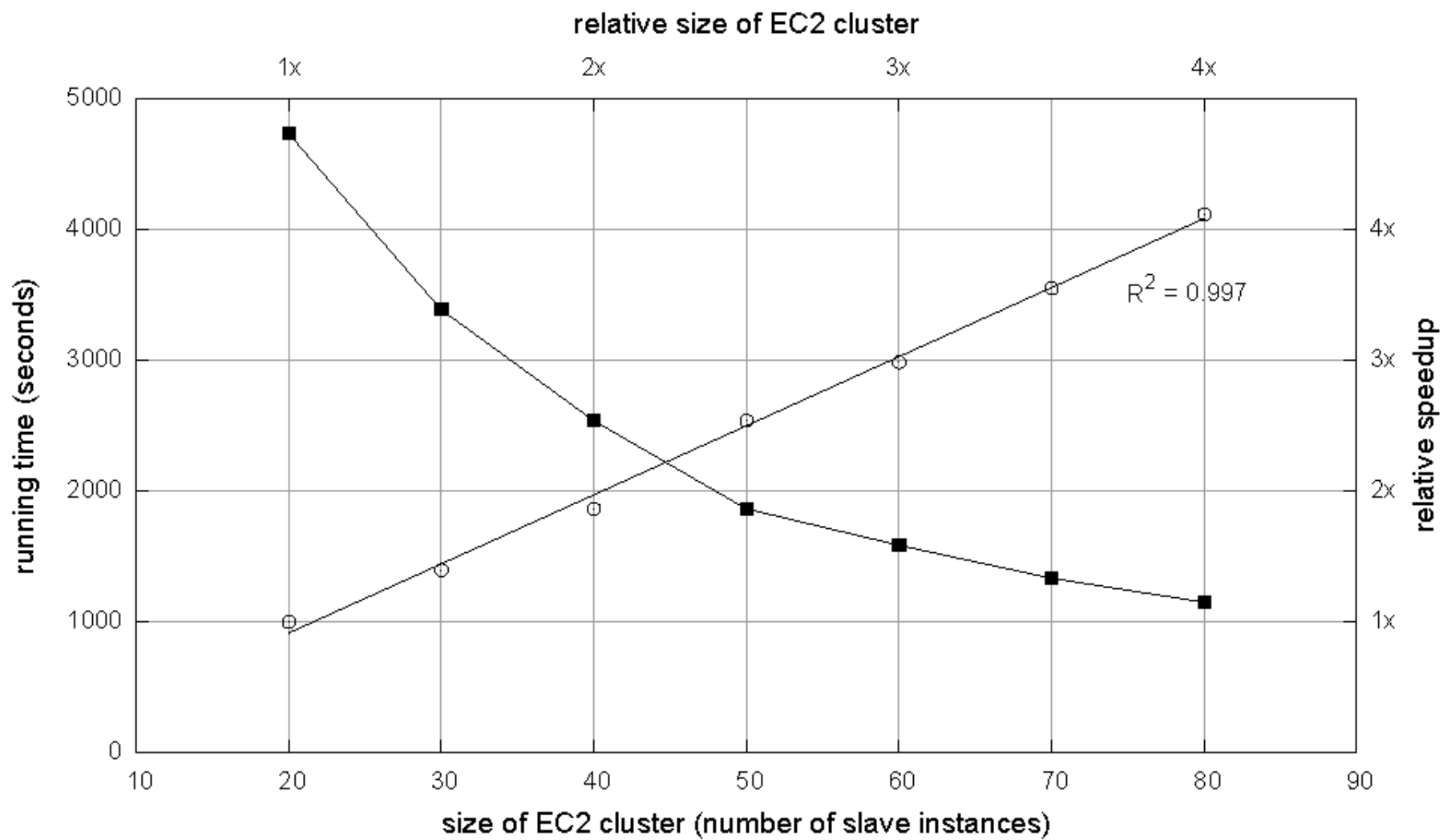


Cluster size: 38 cores

Data Source: Associated Press Worldstream (APW) of the English Gigaword Corpus (v3), which contains 2.27 million documents (1.8 GB compressed, 5.7 GB uncompressed)



## Effect of cluster size on "stripes" algorithm

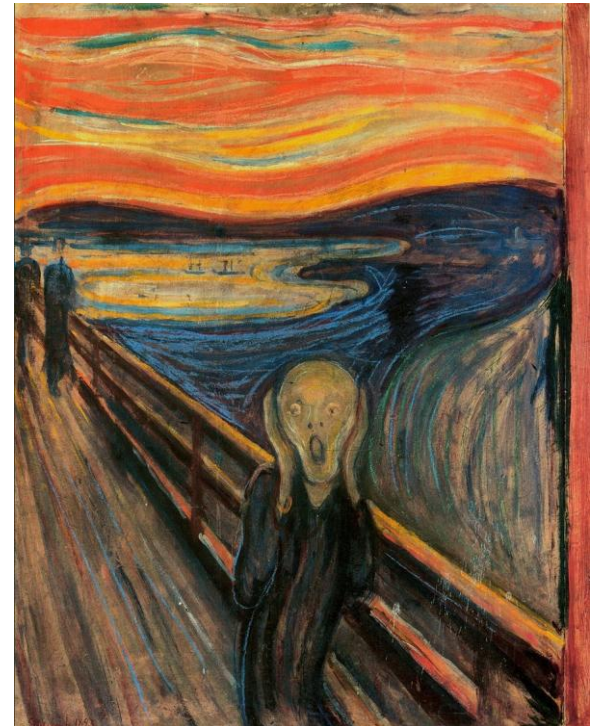


# Debugging at scale

- Works on small datasets, won't scale... why?
  - Memory management issues (buffering and object creation)
  - Too much intermediate data
  - Mangled input records
- Real-world data is messy!
  - There's no such thing as consistent data
  - Watch out for corner cases
  - Isolate unexpected behavior, bring local

# Caveats

- This is bleeding-edge technology (codeword for immature)
  - We have come a long way since 2007, but still far to go
  - Bugs, undocumented “features”, inexplicable behavior, data loss(!)
  - You will experience all these (those W\$\*#T@F! moments)
  - When this happens (and it will)
    - Do not get frustrated (take a deep breath)
    - It’s not the end of the world
- Be patient
  - On a long enough timeline everything works
- Be flexible
  - We will have to be creative in workarounds
- Be constructive
  - Tell me how we can make everyone’s experience better



# Summary

- Further delved into computing using MapReduce
- Introduced map-side optimisations
- Discussed why certain things may not work as expected
- Need to be really careful when designing algorithms to deploy over large datasets
- What seems to work on paper may not be correct when distribution/parallelisation kick in