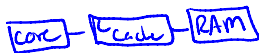


# Midterm

## - Computer Architecture

- CPU is fastest, faster than memory



## - Joins

- Joins are expensive - quadratic input comparisons
- SQL database optimizes join sequence for us
- Filters first then join

## - Maps

maps/dictionaries:  $\{ \text{key: value}, \text{key}_2: \text{value}_2 \}$

Tuples (unchangeable) can have keys

map is much faster than  
 $|S| + |T|$

Joining more than one field e.g. (first + last name)

- Cartesian  $n \times m$



- Big Data requires reading 2 blocks at a time

- Making algorithms more efficient:

- order of evaluation
- minimize intermediate results
- SQL database does the best sequence for user
- pandas we need to control it on our own

\* Joins iterate over both tables

- Index makes lookup more efficient if join is on index key

## - Large Volumes of data

1. Reduce amount of data
  2. Use more efficient algos or data structures
  3. Split work among multiple processors
- Modern computers have 4 CPU cores

## - Parallelism - Multicore Processing

- apply function on column uses parallelism by running function on different rows on different cores
- results then merged together
- DASH parallel PANDAS
  - partitions data
  - partition (e.g. 100) then merge
  - partitions needs to be more than # of cores
- Google Colab runs on multicore machine
- Pandas in itself doesn't use multicore

## - Relational Algebra

- Groupby
- join
- select
- apply
- merge