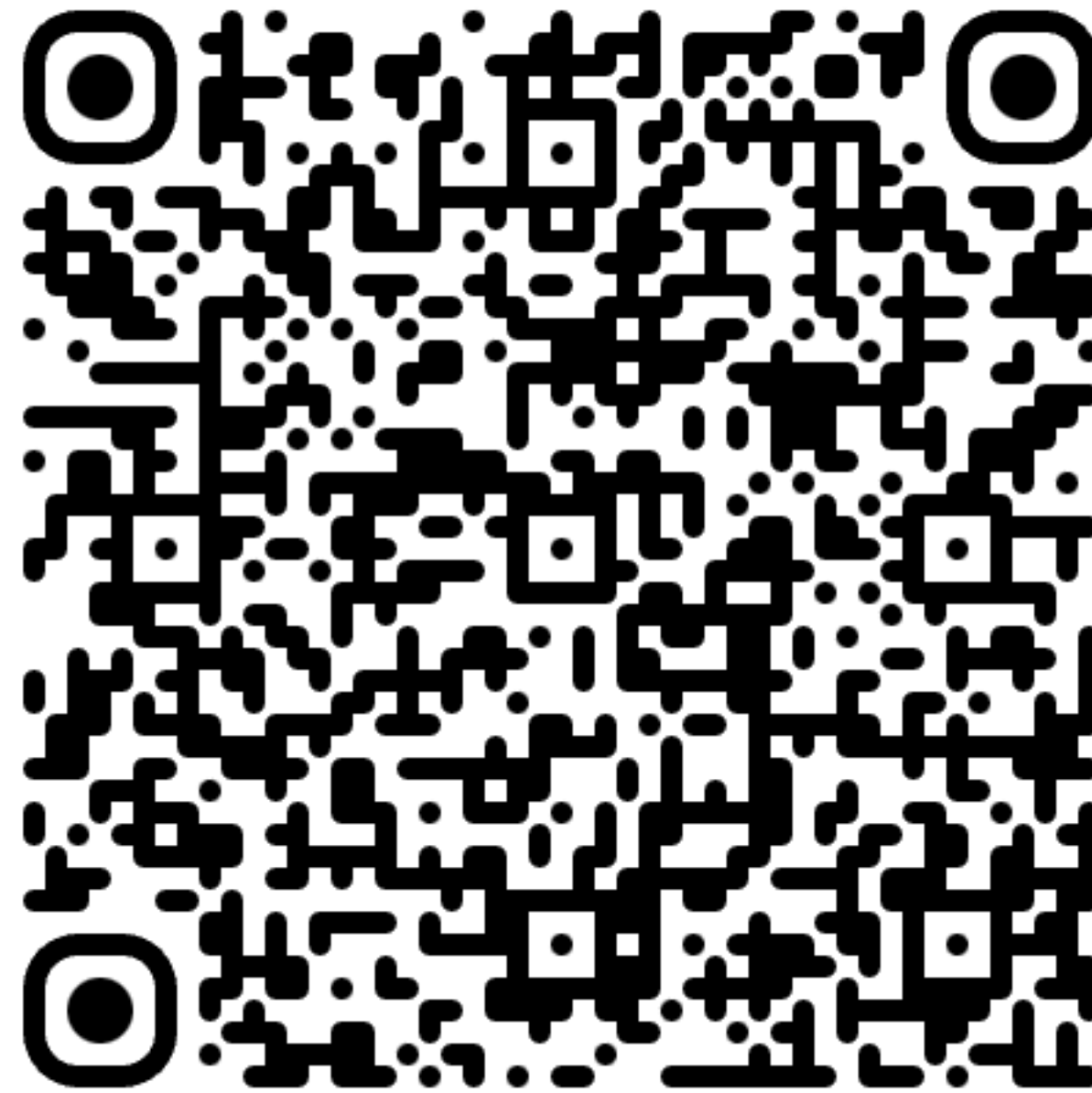


Coding for Economists

Advanced Session 4

Jian Cao
19 May 2025

Module Files



[Google Drive Folder](#)

Outline

Outline

- Large Data Problems

Outline

- Large Data Problems
- Efficient Format / Operation

Outline

- Large Data Problems
- Efficient Format / Operation
- Efficient I/O

Outline

- Large Data Problems
- Efficient Format / Operation
- Efficient I/O
- Parallel / Distributed Computing

Large Data Problems

Large Data Problems

- **Memory**-Bound

Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory

Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound

Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

Large Data Problems

Solution

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

Solution

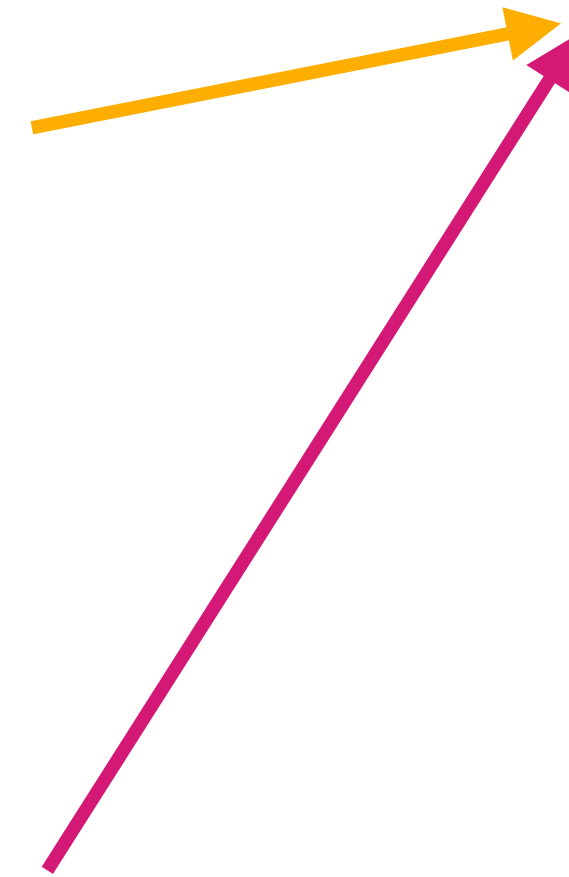
- Pandas, Numpy format

Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

Solution

- Pandas, Numpy format

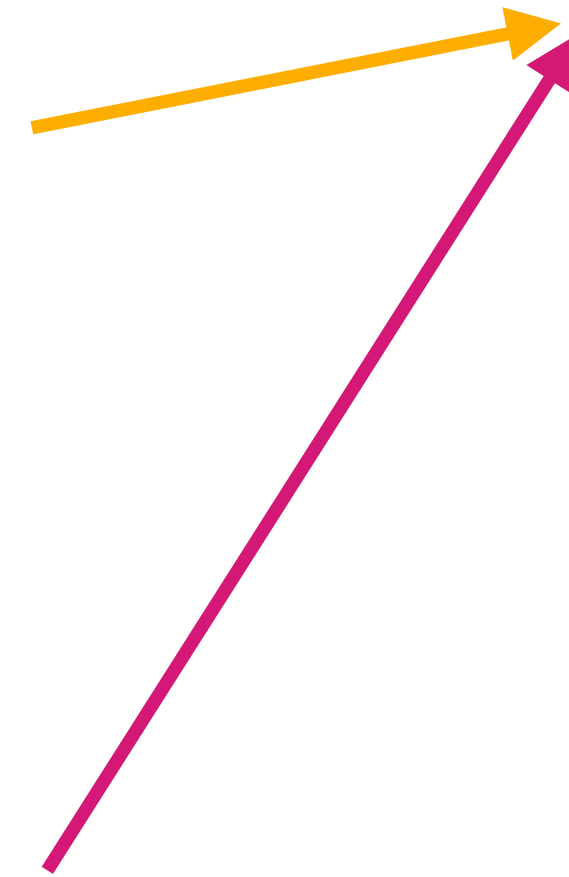


Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

Solution

- Pandas, Numpy format
- Vectorized operation

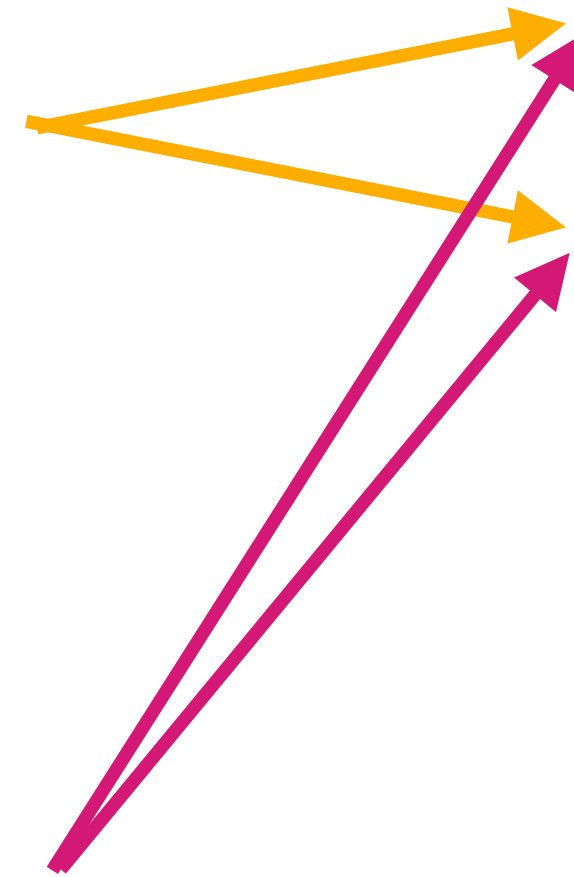


Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

Solution

- Pandas, Numpy format
- Vectorized operation

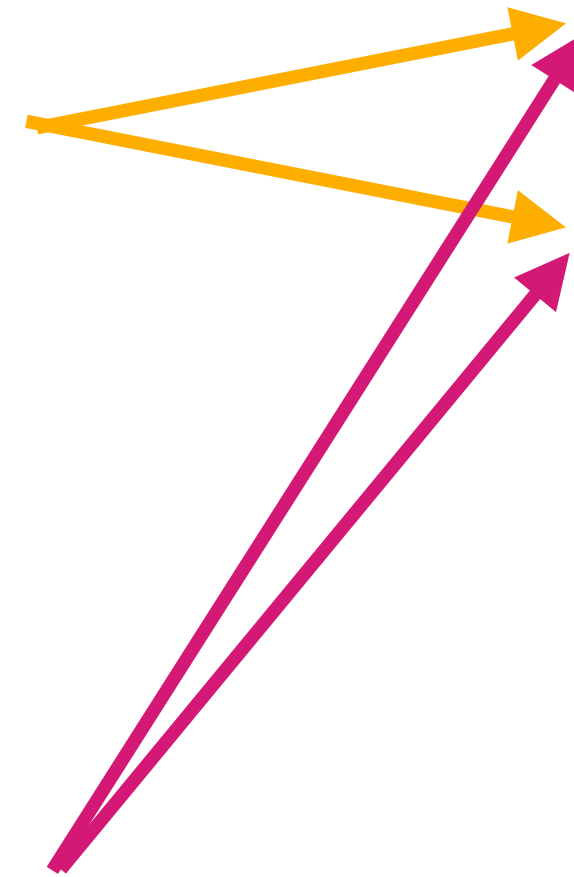


Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

Solution

- Pandas, Numpy format
- Vectorized operation
- Chunked I/O

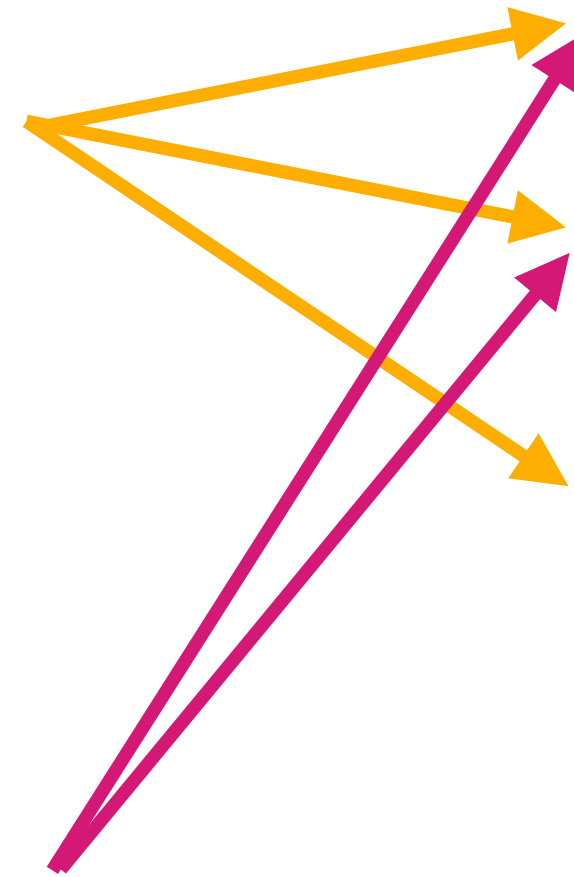


Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

Solution

- Pandas, Numpy format
- Vectorized operation
- Chunked I/O

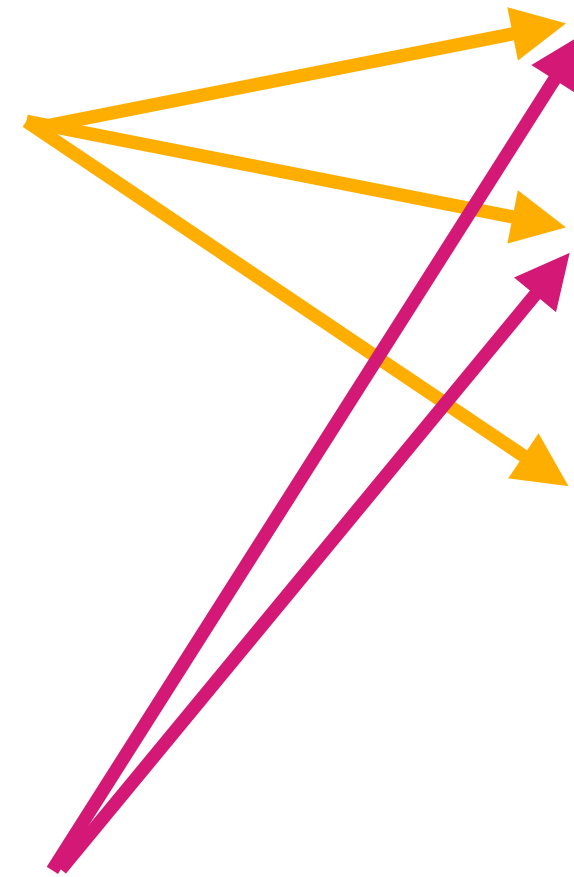


Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

Solution

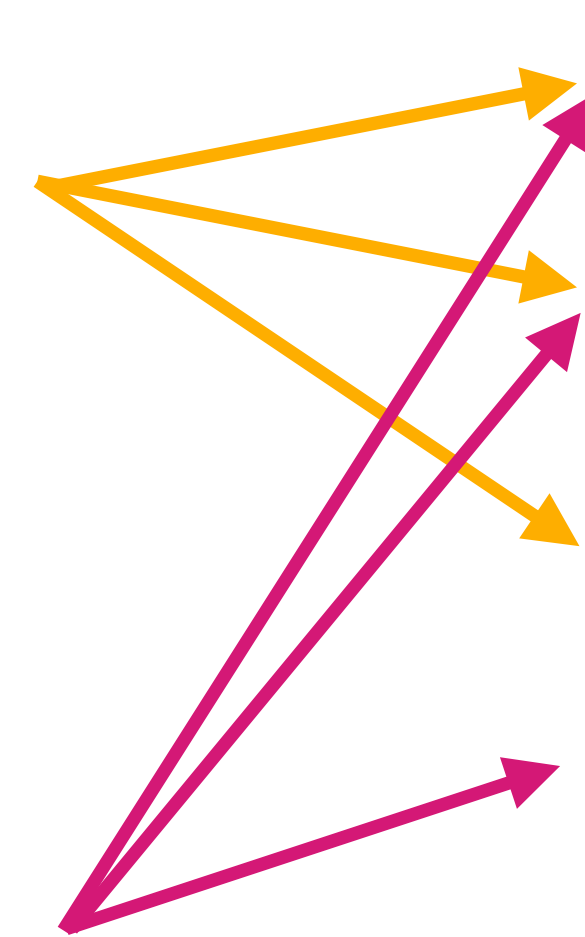
- Pandas, Numpy format
- Vectorized operation
- Chunked I/O
- Parallel computing



Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

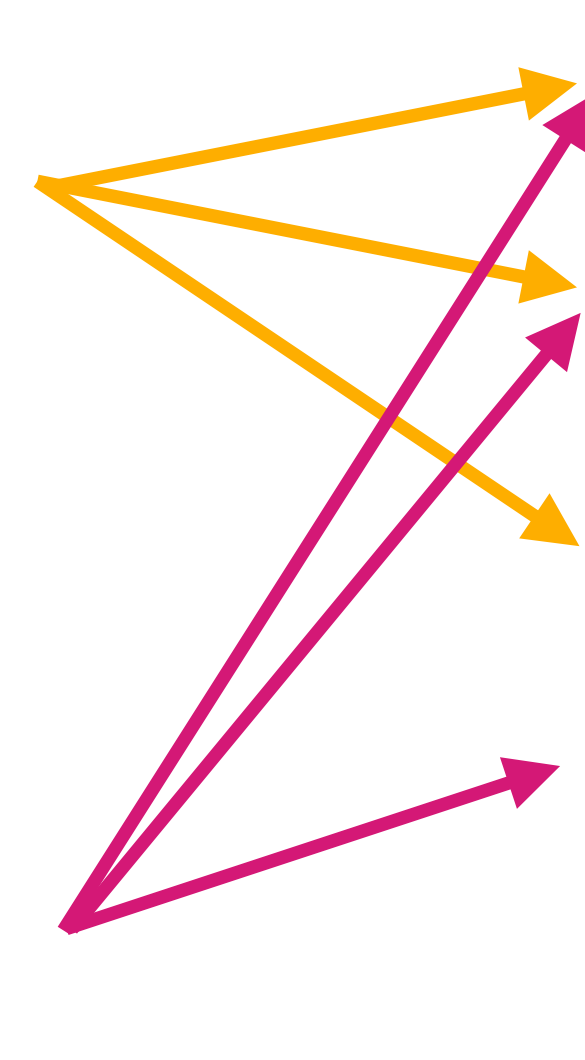
Solution

- 
- The diagram consists of two sets of arrows. Three yellow arrows originate from the 'Memory-Bound' section and point to the first three solutions: 'Pandas, Numpy format', 'Vectorized operation', and 'Chunked I/O'. Two magenta arrows originate from the 'Compute-Bound' section and point to the last two solutions: 'Parallel computing' and 'Vectorized operation'.
- Pandas, Numpy format
 - Vectorized operation
 - Chunked I/O
 - Parallel computing

Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

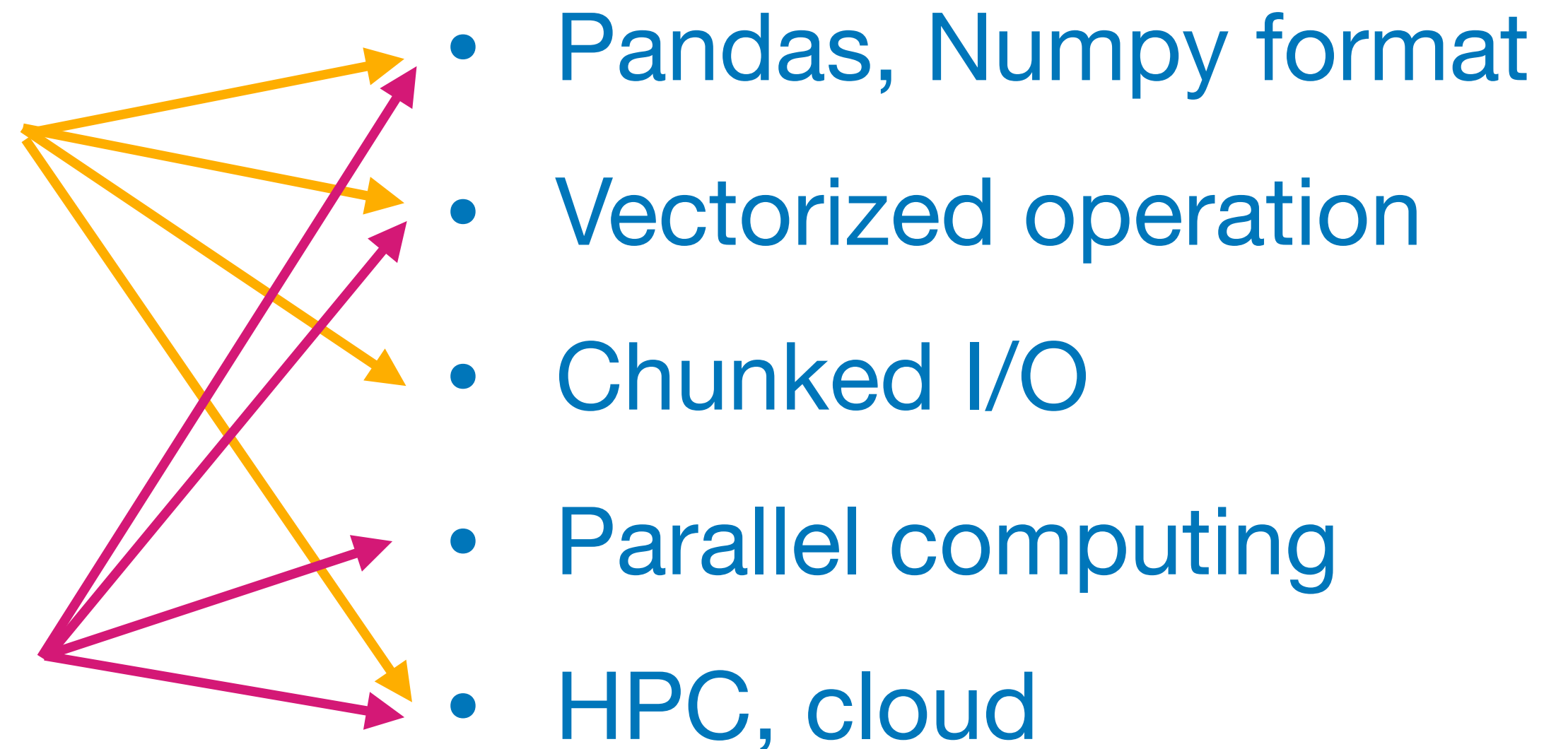
Solution

- 
- A diagram consisting of two sets of arrows. Three yellow arrows originate from the 'Memory-Bound' problem and point to the first three solutions: 'Pandas, Numpy format', 'Vectorized operation', and 'Chunked I/O'. Two magenta arrows originate from the 'Compute-Bound' problem and point to the last two solutions: 'Parallel computing' and 'HPC, cloud'.
- Pandas, Numpy format
 - Vectorized operation
 - Chunked I/O
 - Parallel computing
 - HPC, cloud

Large Data Problems

- **Memory**-Bound
 - Program takes long time moving files in/out RAM, or crashes due to out of memory
- **Compute**-Bound
 - Program takes long time waiting for computing cores

Solution



Efficient Format / Operation

Efficient Format / Operation

- Pandas, Numpy data format

Efficient Format / Operation

- Pandas, Numpy data format
 - Homogeneous, contiguous

Efficient Format / Operation

- Pandas, Numpy data format
 - Homogeneous, contiguous
 - Faster computation

Efficient Format / Operation

- Pandas, Numpy data format
 - Homogeneous, contiguous
 - Faster computation
- Vectorized operation

Efficient Format / Operation

- Pandas, Numpy data format
 - Homogeneous, contiguous
 - Faster computation
- Vectorized operation
 - Vector, matrix >> element loops

Efficient Format / Operation

- Pandas, Numpy data format
 - Homogeneous, contiguous
 - Faster computation
- Vectorized operation
 - Vector, matrix >> element loops
 - Calls C++, CUDA

Efficient Format / Operation

- Pandas, Numpy data format
 - Homogeneous, contiguous
 - Faster computation
- Vectorized operation
 - Vector, matrix >> element loops
 - Calls C++, CUDA
 - Bulk memory access

Efficient Format / Operation

- Pandas, Numpy data format
 - Homogeneous, contiguous
 - Faster computation
- Vectorized operation
 - Vector, matrix >> element loops
 - Calls C++, CUDA
 - Bulk memory access

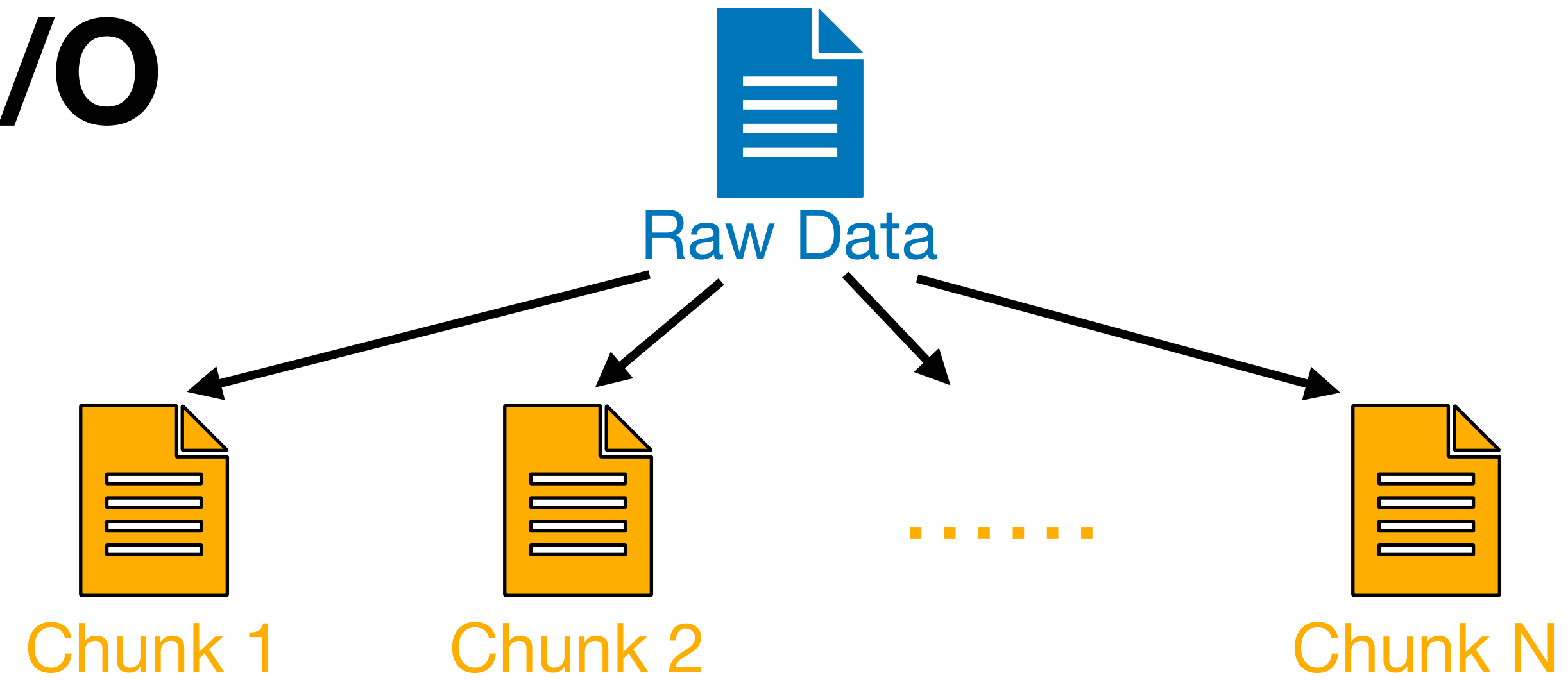
Chunked I/O

Chunked I/O

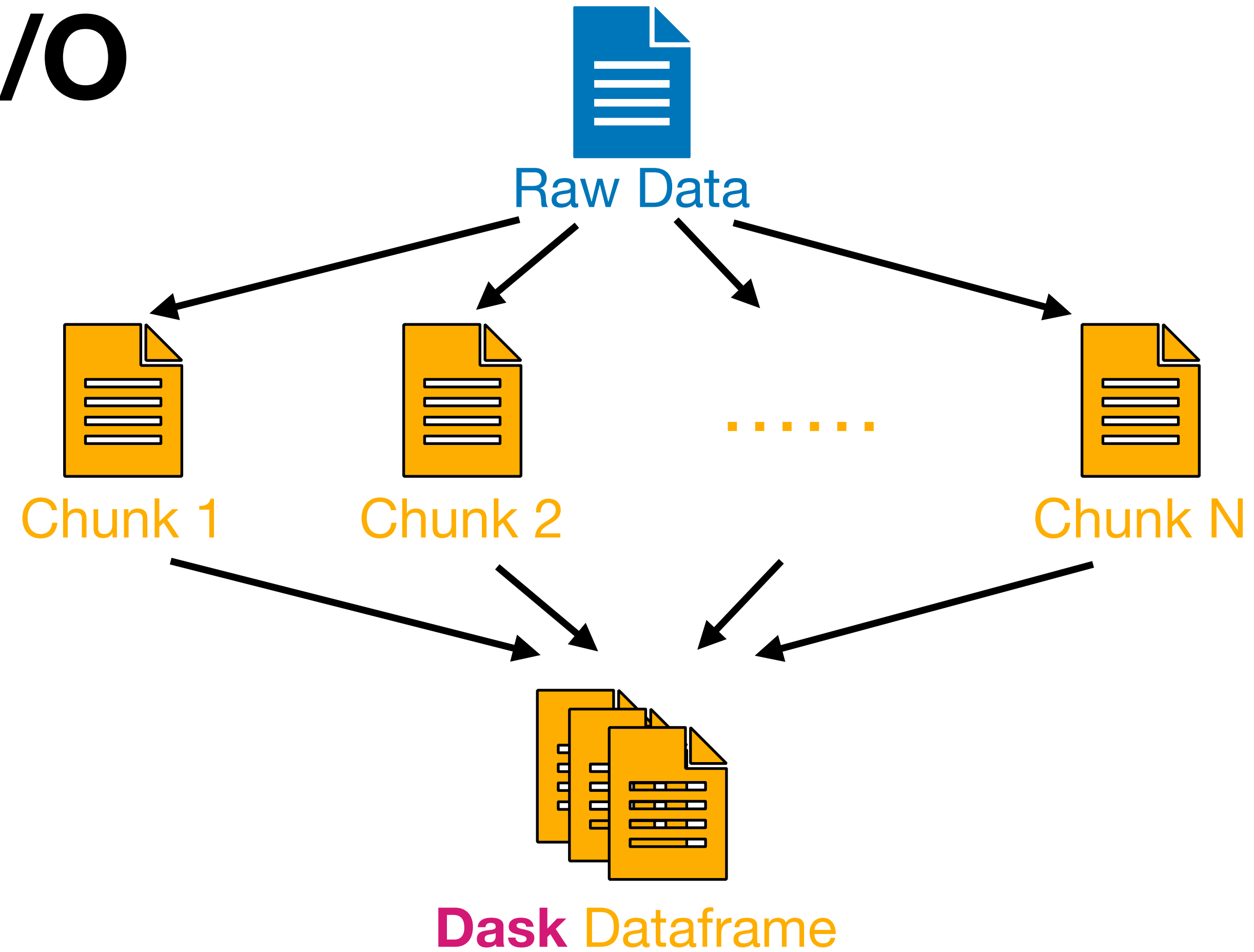


Raw Data

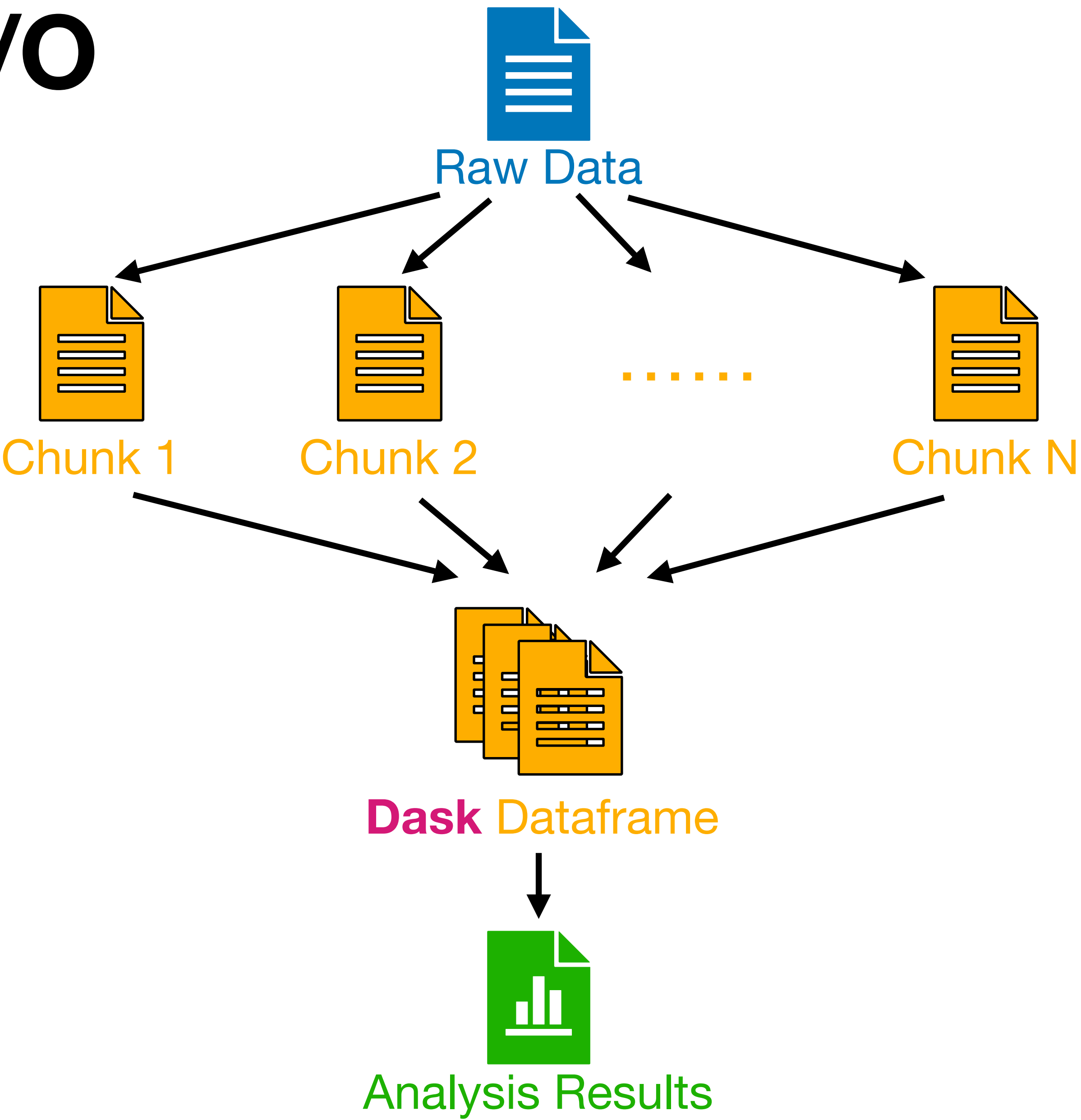
Chunked I/O



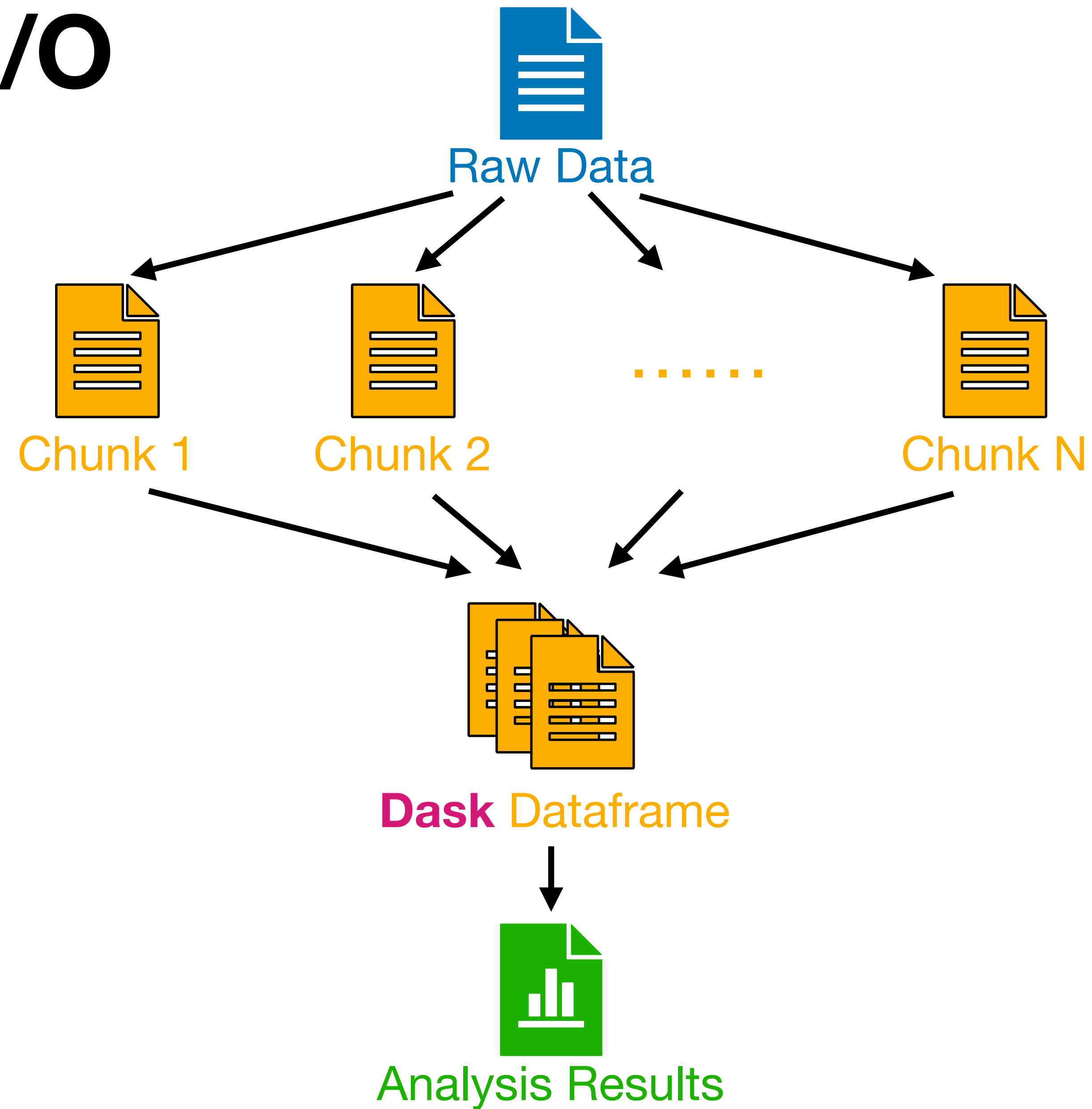
Chunked I/O



Chunked I/O



Chunked I/O

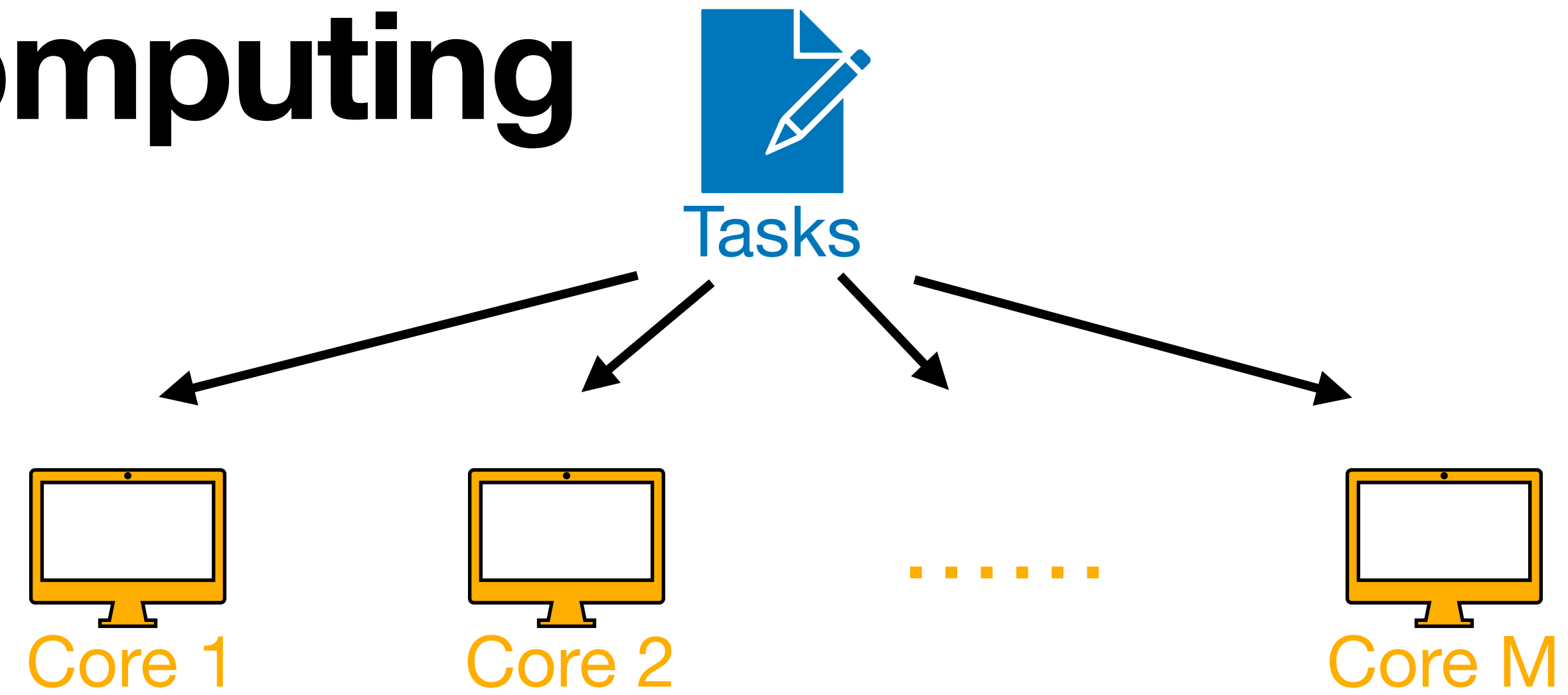


Parallel Computing

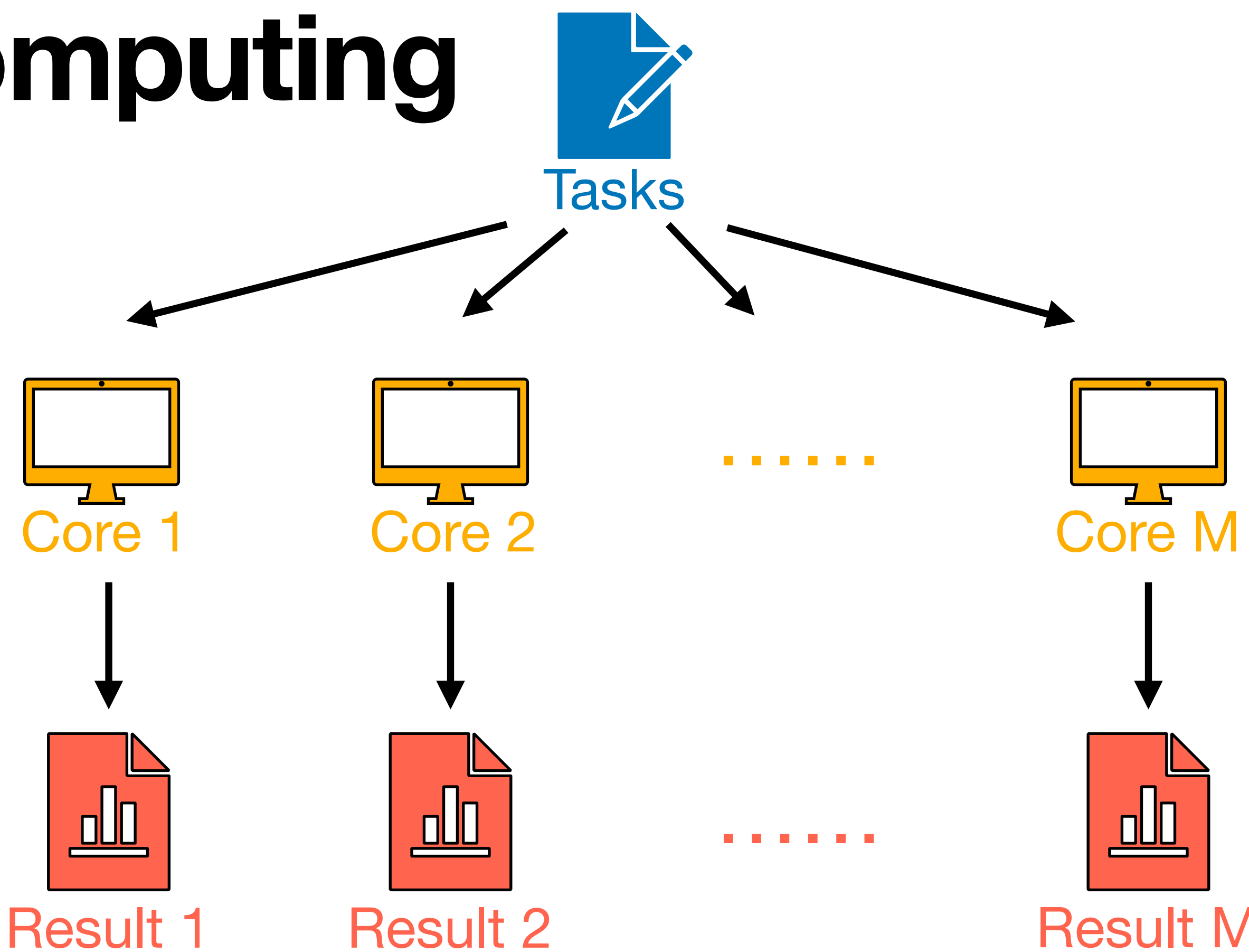
Parallel Computing



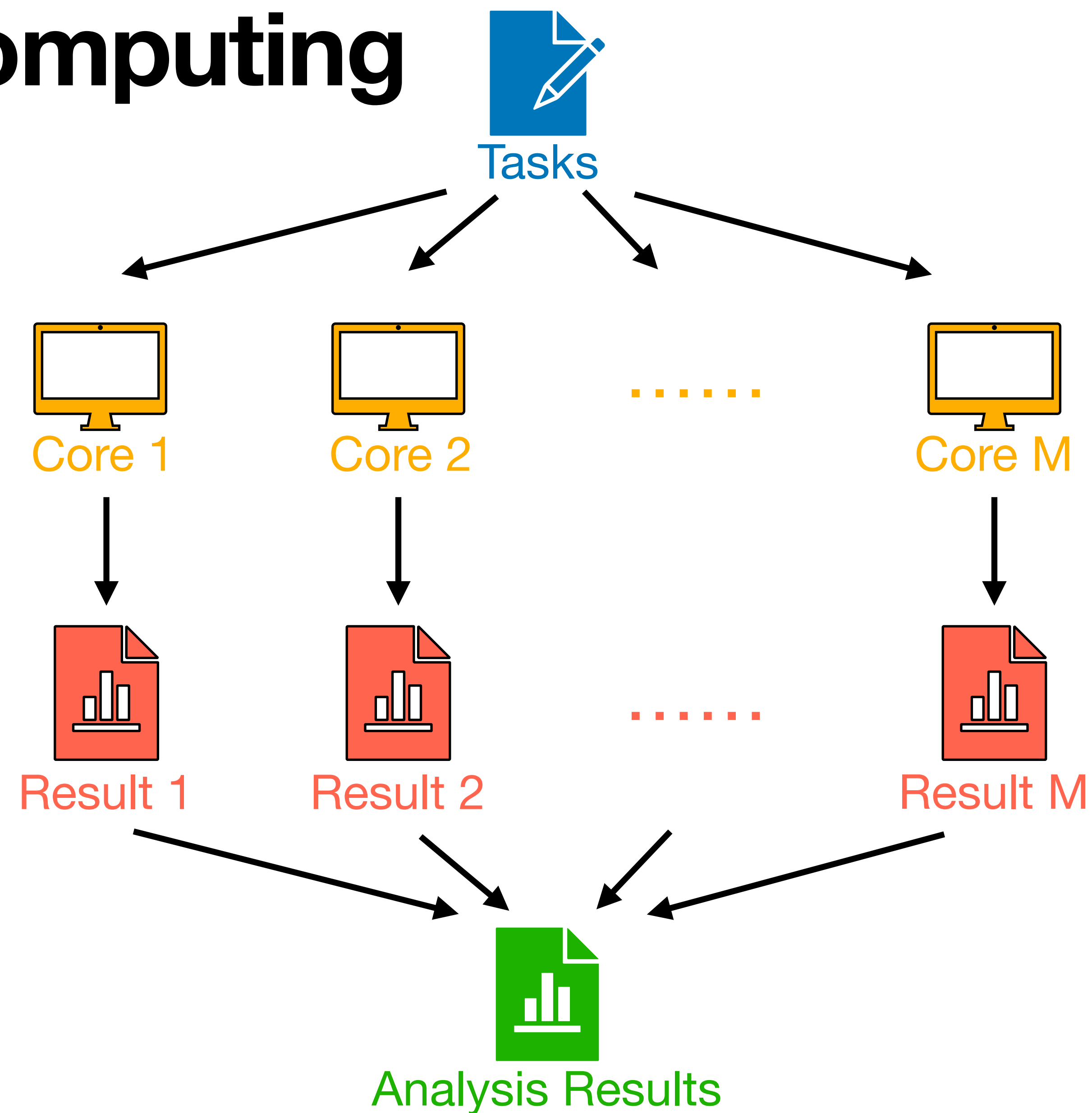
Parallel Computing



Parallel Computing



Parallel Computing



Parallel Computing

