# Metagenome is the genome of a microbial community

# Metagenome sequencing



Harvest microbes

Microbes

Extract DNA

Genomic DNA

Shear, & Sequencing

Short Reads

Assembly

Reconstructed genomes

# Metagenome assembly

Genome ~= Book   Metagenome ~= Library

Library of Books        Shredded Library        "reconstructed" Library
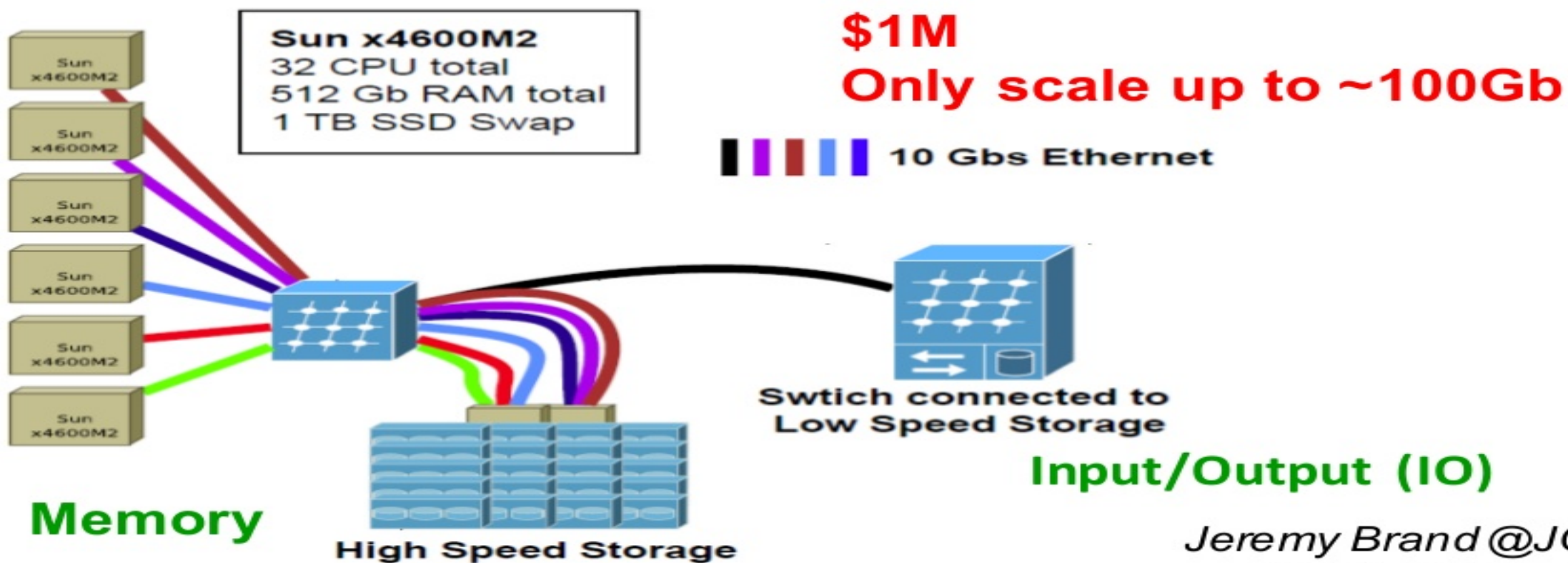


Sequencing ~= sampling the pieces

# Complexity is another…

- Data complexity
  - Contamination
  - Number of microbial species
  - Species abundance distribution
  - Sequencing errors
- Algorithm complexity
  - Multiple steps, each has different time/space characteristics

# The Ideal Solution

- ☐ Easy to develop
- ☐ Robust
- ☐ Scale to big data
- ☐ Efficient

# 2009: Special Hardware

**Sun x4600M2**
32 CPU total
512 Gb RAM total
1 TB SSD Swap

**$1M**
**Only scale up to ~100Gb**

| | | | | | 10 Gbs Ethernet

**Memory**

**High Speed Storage**

**Swtich connected to Low Speed Storage**

**Input/Output (IO)**

*Jeremy Brand @JGI*
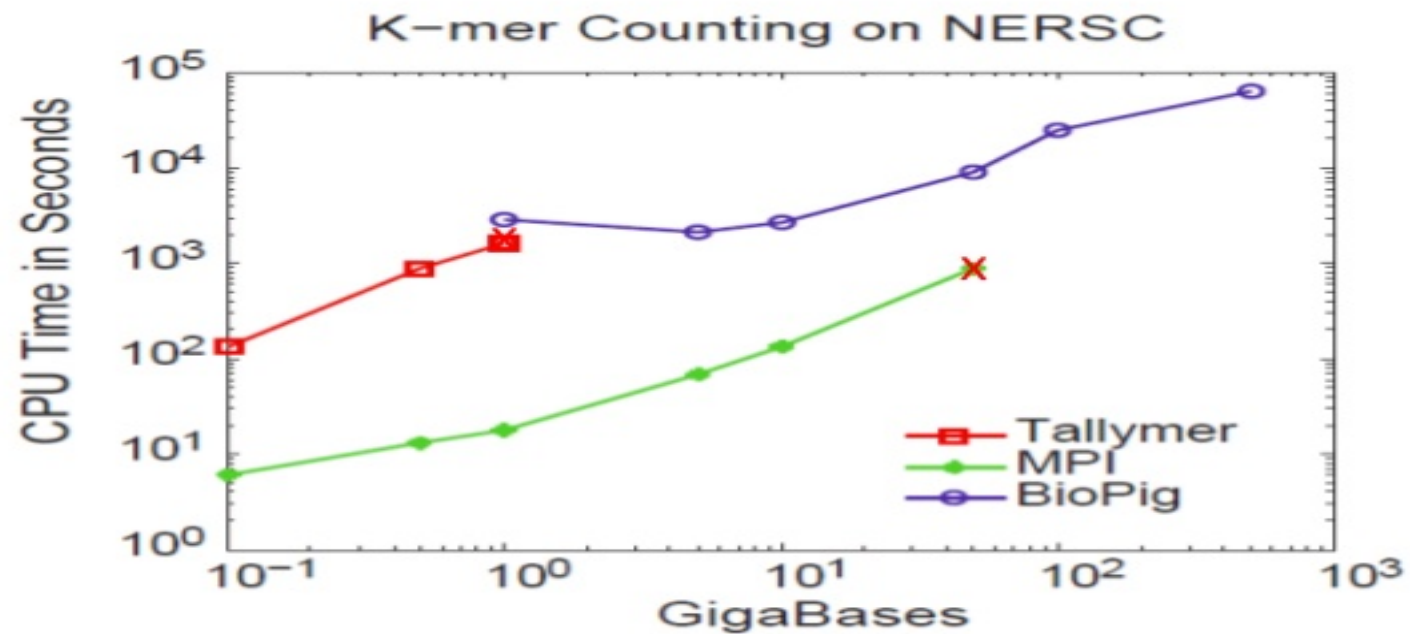*FPGA @Convey*

# 2010: MP/MPI on supercomputers



**MPI version
412 Gb, 4.5B reads
2.7 hours on 128x24 cores
NESRC Supercomputer**

*Fast, scalable*

Problems:

- **Experienced software engineers**
- **Six months of development time**
- **One task fails, all tasks fail**

*Rob Egan @JGI*

# 2013: Hadoop/BioPig



K−mer Counting on NERSC

*Karan Bhatia, Henrik Nordberg, Kai Wang*

# Challenges in application

- <span style="color:red">2-3 orders of magnitude slower than MPI</span>
- IO optimization, e.g., reduce data copying
- Some problems do not easily fit into map/reduce framework, e.g., graph-based algorithms
- Runs on AWS, but cost $$$ if not optimized

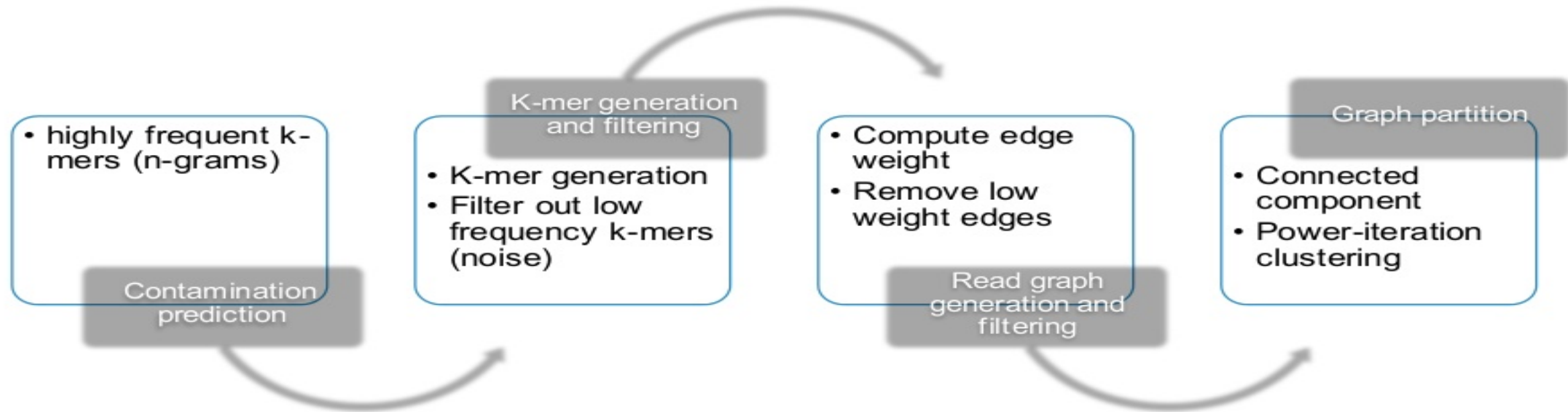# Addressing big data: Apache Spark

- New **scalable** programming paradigm
  - Compatible with Hadoop-supported storage systems
- Improves **efficiency** through:
  - In-memory computing primitives
  - General computation graphs
- Improves **usability** through:
  - Rich APIs in Java, Scala, Python
  - Interactive shell

☐ **Scale to big data**

☐ **Efficient**

☐ **Easy to develop**

☐ **Robust?**

# Goal: Metagenome read clustering

- Clustering reads based on their genome of origin can reduce metagenome problem to single-genome problem
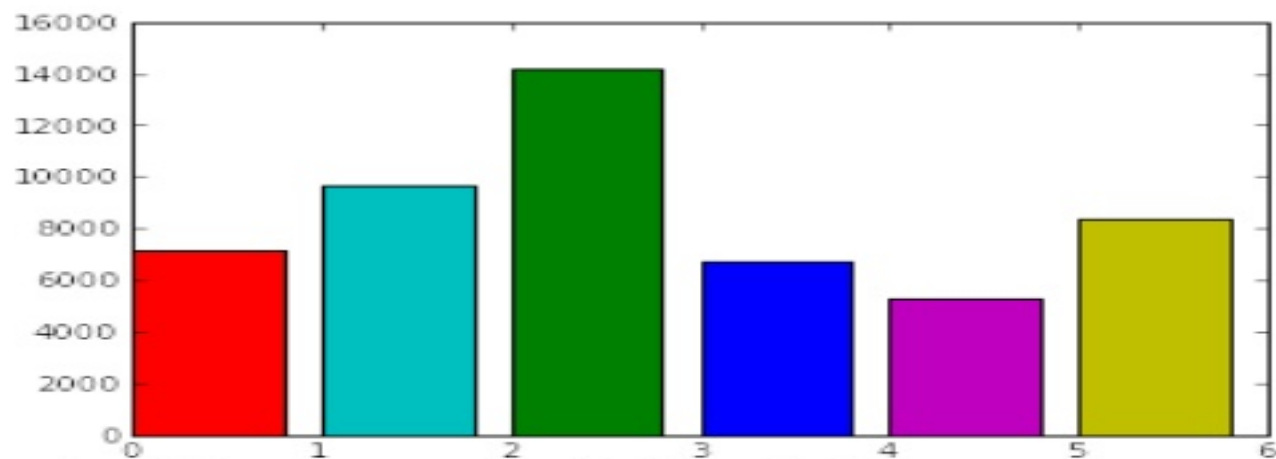- Ideally scale up to TB data sizes

# Algorithm

- highly frequent k-mers (n-grams)

**Contamination prediction**

**K-mer generation and filtering**
- K-mer generation
- Filter out low frequency k-mers (noise)

- Compute edge weight
- Remove low weight edges

**Read graph generation and filtering**

**Graph partition**
- Connected component
- Power-iteration clustering

# Platforms we run spark

- Standalone Spark on single large memory server
- On-demand Spark cluster over HPC
- AWS Elastic Map Reduce (EMR)

# Testing the accuracy of the algorithm with a small toy test dataset

- Species:
  - 6 bacterial species (10kb from each)
  - Synthetic communities with random proportions of each genome, reads drawn from single genome sequencing projects (noisy)
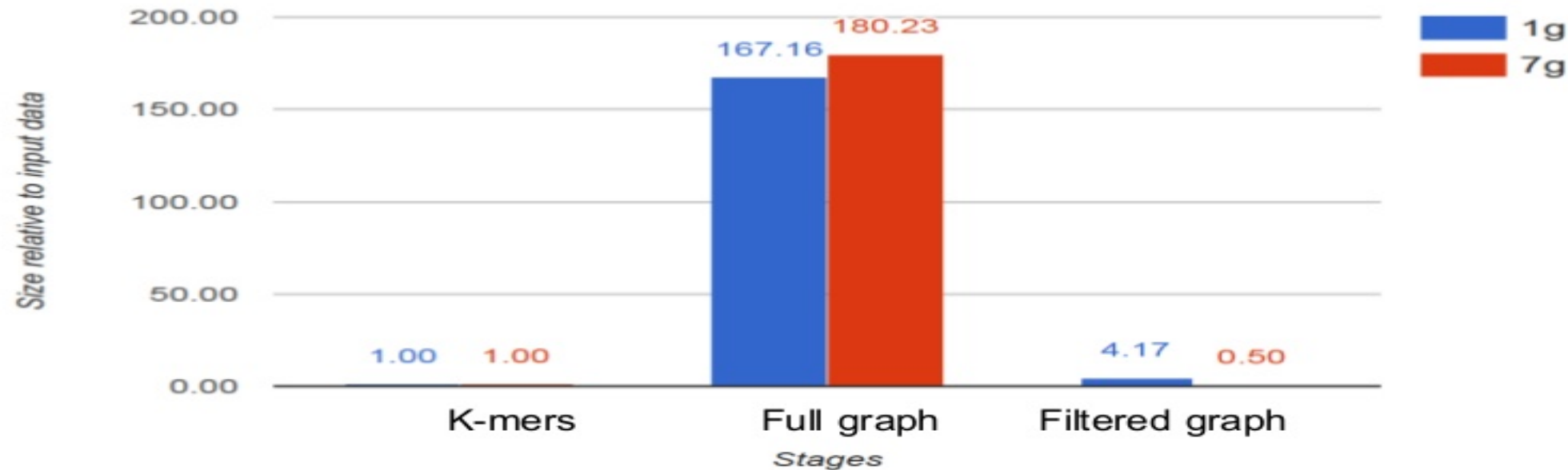  - Ideal situation (no shared sequences between genomes, sufficient sequencing coverage):



Reads of the same color belong to the same genome

# Real world datasets

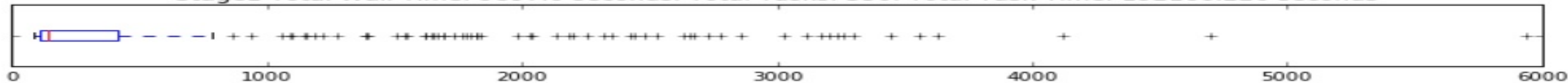| Dataset | Number of spieces | Sampling depth |
|---|---|---|
| Soil metagenome | High | Low |
| Cow rumen metagenome | Medium | Low to medium |
| Maize transcriptome ("fake metagenome") | Low | High |

# Tuning parallelism for load balance

## K-mer generation/filtering: partition_size = 16Mb



Stage1 Total Wall Time: 9897.0 seconds. Total Tasks: 350. Total Task Time: 192166.226 seconds
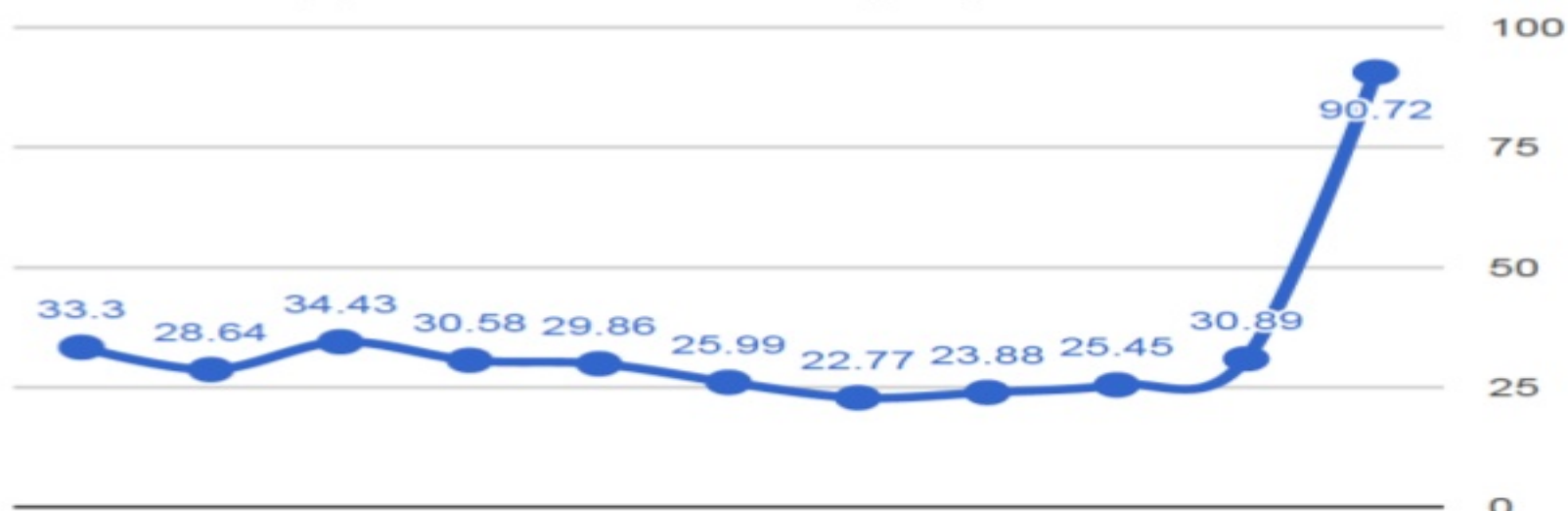
## K-mer generation/filtering: partition_size = 1.3Mb



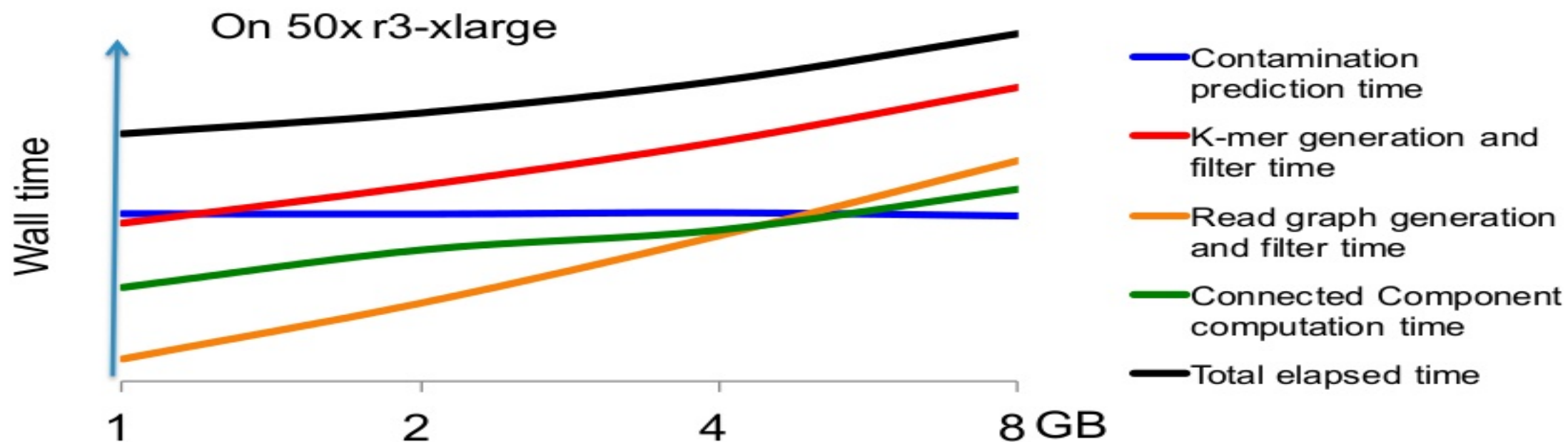Stage9 Total Wall Time: 554.0 seconds. Total Tasks: 4200. Total Task Time: 152192.266 seconds

# Tuning parallelism for performance

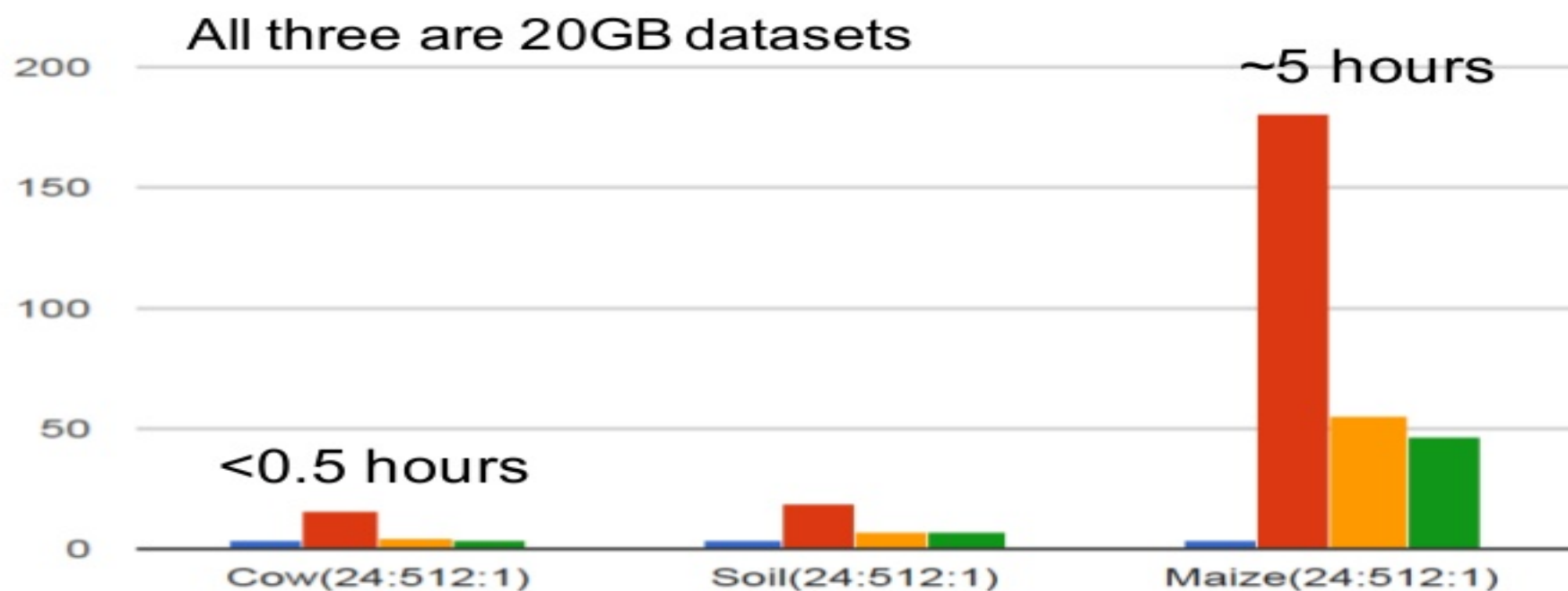Decreasing partition size for graph construction



Optimizing parameters can reduce total running time from > 90 min to 20 min
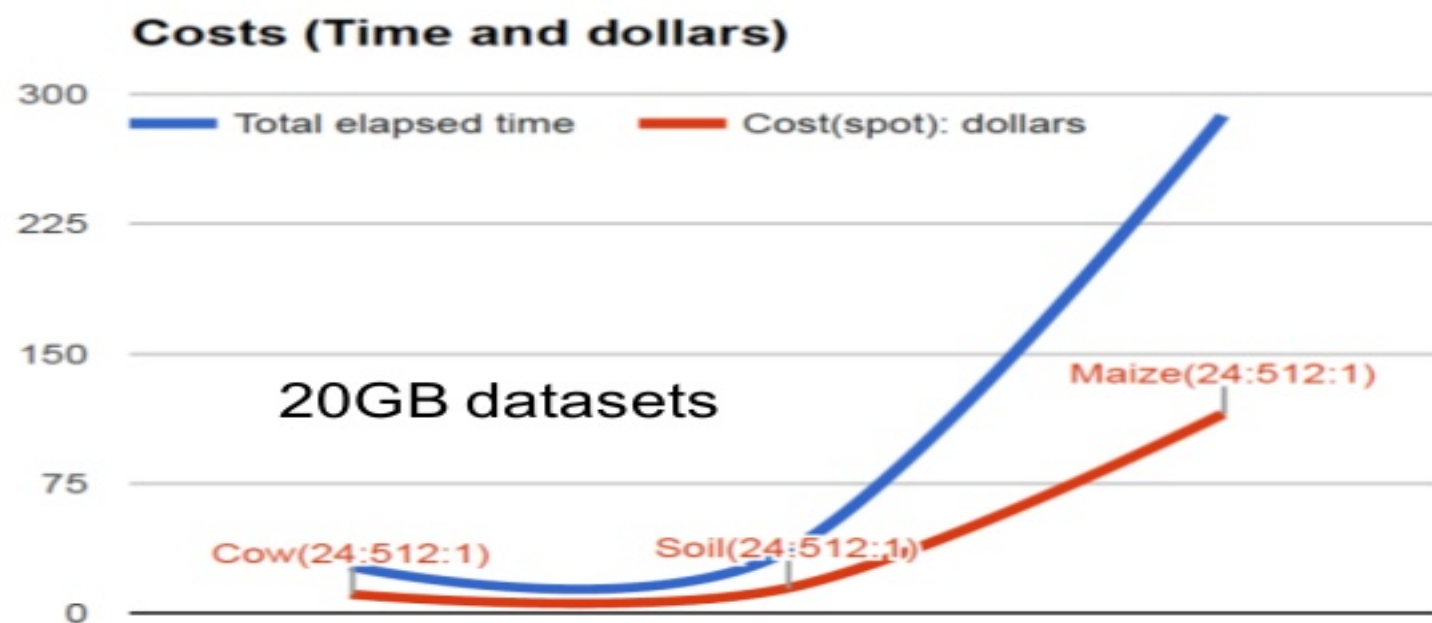
# Scale well on small data



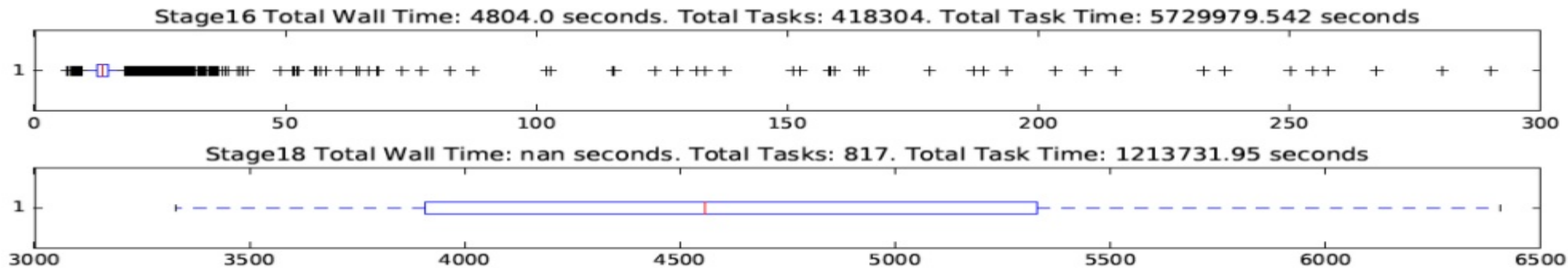On 50x r3-xlarge

# Cost on AWS EMR spot instances



Costs (Time and dollars)

20GB datasets

Projected cost for a typical 1TB dataset: ~$500

# Scaling up to 100Gb: failed

Read graph generation/filtering



Stage16 Total Wall Time: 4804.0 seconds. Total Tasks: 418304. Total Task Time: 5729979.542 seconds

Stage18 Total Wall Time: nan seconds. Total Tasks: 817. Total Task Time: 1213731.95 seconds

Reducing partition numbers for graph partitioning

# Potential solutions

- Avoid shuffling:
  - generate the graph, save to disk, then merge partitions outside of Spark.
- Size-specific parameters
  - Larger datasets may not use parallelism parameters optimized for smaller ones
- Your inputs…

# Acknowledgements

**Spark Team**
Xiandong Meng
Jordan Hoffman@Harvard

Lisa Gerhardt , Evan Racah
@ NERSC

Gary Jung, Greg Kurtzer
Bernard Li,Yong Qin @ HPC