# The Impact of HPC Best Practice Applied to Next-Generation Sequencing Workflows

- <u>Pierre Carrier</u>, Bill Long, Richard Walsh, Jef Dawson

- Carlos P. Sosa

- Brian Haas, Timothy Tickle

- Thomas William

RNA-Seq software:

# Agenda

- What have we **learned in computational chemistry**?
- A tale of **two communities**
- Why **HPC Best Practice**?

http://biorxiv.org/content/early/2015/04/07/017665

# Agenda

- What have we **learned in computational chemistry**?
- A tale of **two communities**
- Why **HPC Best Practice**?
- What are **NGS workflows**?
- What is **Trinity RNA-Seq**?

http://biorxiv.org/content/early/2015/04/07/017665

# Agenda

- What have we **learned in computational chemistry**?
- A tale of **two communities**
- Why **HPC Best Practice**?
- What are **NGS workflows**?
- What is **Trinity RNA-Seq**?
- **HPC Best Practice in a nutshell**
- **MPI-Inchworm** algorithm

bioR𝛘iv
beta

**THE PREPRINT SERVER FOR BIOLOGY**

http://biorxiv.org/content/early/2015/04/07/017665

# Agenda

- What have we **learned in computational chemistry**?
- A tale of **two communities**
- Why **HPC Best Practice**?
- What are **NGS workflows**?
- What is **Trinity RNA-Seq**?
- **HPC Best Practice in a nutshell**
- **MPI-Inchworm** algorithm
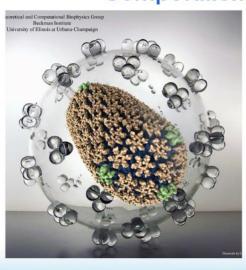- Hardware resources
- **Scaling** memory and CPU
- Conclusions

http://biorxiv.org/content/early/2015/04/07/017665

# Computational Chemistry Embraces HPC

# A Tale of Two Communities

**Computational Chemistry**

**Computational Genomics**

Method implementation has evolved together with HPC

Method implementation has evolved together with sequencing technologies (NGS), not with HPC

# A Tale of Two Communities

**Computational Chemistry**

**Computational Genomics**

Method implementation has evolved together with HPC

⟷

Method implementation has evolved together with sequencing technologies (NGS), not with HPC

Example: Hartree-Fock (HF) and Density Functional Theory (DFT).
The initial problem, HF, is essentially very large to solve, often too large.
DFT, was created (in part) in order to more easily adapt the problem size and accuracy (xc functionals) to the available supercomputer

# A Tale of Two Communities

**Computational Chemistry**

**Computational Genomics**

Method implementation has evolved together with HPC

⟷

Method implementation has evolved together with sequencing technologies (NGS), not with HPC

# A Tale of Two Communities

## Computational Chemistry

## Computational Genomics

Method implementation has evolved together with HPC

$\longleftrightarrow$

Method implementation has evolved together with sequencing technologies (NGS), not with HPC

2 main "next-generation" sequencing technologies:

Short reads

Long reads

A "Reads" = a "puzzle piece" of nucleotides (the "ATGC"s)

illumina®

PACIFIC BIOSCIENCES®
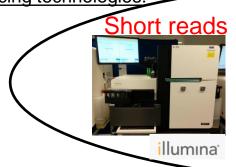
# A Tale of Two Communities

**Computational Chemistry**

**Computational Genomics**

Method implementation has evolved together with HPC

⟷

Method implementation has evolved together with sequencing technologies (NGS), not with HPC

2 main "next-generation" sequencing technologies:

**Short reads**

**Long reads**

A "Reads" = a "puzzle piece" of nucleotides (the "ATGC"s)

Short reads: ~100s nucleotides

Long reads: ~1000s nucleotides

illumina

PACIFIC BIOSCIENCES

# A Tale of Two Communities

**Computational Chemistry**

**Computational Genomics**

Method implementation has evolved together with HPC

⟷

Method implementation has evolved together with <u>sequencing technologies</u> (NGS), not with HPC

<u>Emerging technologies:</u>

Short reads

A "<u>Reads</u>" = a "puzzle piece"
of nucleotides (the "ATGC"s)

Long reads

…Even longer reads

Short reads:
~**100s** nucleotides

Long reads:
~**1000s** nucleotides

illumina

PACIFIC BIOSCIENCES

NANOPORE

# A Tale of Two Communities

## Computational Chemistry

## Computational Genomics

Method implementation has evolved together with HPC

←→

Method implementation has evolved together with sequencing technologies (NGS), not with HPC

**Short reads**

**sequencing**
- ✓ Smallest genomes
- ✓ Bacteria
- ✓ Drosophila
- ✓ Mouse
- ✓ **Human**
- ✓ …
- ✓ Spruce
- X Axolotl (not done)
- X Largest genome

SIZE ↓

**Software**
- ✓ Develop/compute on laptop
- ✓ On laptop
- ✓ On laptop
- ✓ On workstation
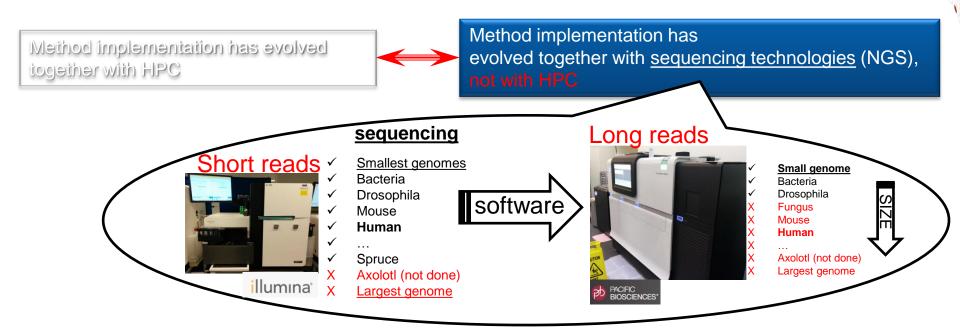- ✓ **On workstation**
- ✓ …
- ✓ On supercomputer?

illumina®

# A Tale of Two Communities

**Computational Chemistry**

**Computational Genomics**

Method implementation has evolved together with HPC

⬌

Method implementation has evolved together with sequencing technologies (NGS), not with HPC

**sequencing**

Short reads

✓ Smallest genomes
✓ Bacteria
✓ Drosophila
✓ Mouse
✓ **Human**
✓ …
✓ Spruce
X Axolotl (not done)
X Largest genome

illumina

software →

Long reads

✓ **Small genome**
✓ Bacteria
✓ Drosophila
X Fungus
X Mouse
X **Human**
X …
X Axolotl (not done)
X Largest genome

SIZE

PACIFIC BIOSCIENCES

# A Tale of Two Communities

## Computational Chemistry

## Computational Genomics

Method implementation has evolved together with HPC

⟷

Method implementation has evolved together with sequencing technologies (NGS), not with HPC

HPC best practice is a priority when developing software

⟷

Sequencing technologies are rapidly and continually changing. It forces developers to focus on functionality rather than system size (i.e., HPC best practice)

# A Tale of Two Communities

**Computational Chemistry**

**Computational Genomics**

Method implementation has evolved together with HPC

↔

Method implementation has evolved together with <u>sequencing technologies</u> (NGS), not with HPC

HPC best practice is a priority when developing software

↔

Sequencing technologies are rapidly and continually changing. It forces developers to focus on <u>functionality</u> rather than system size (i.e., HPC best practice)

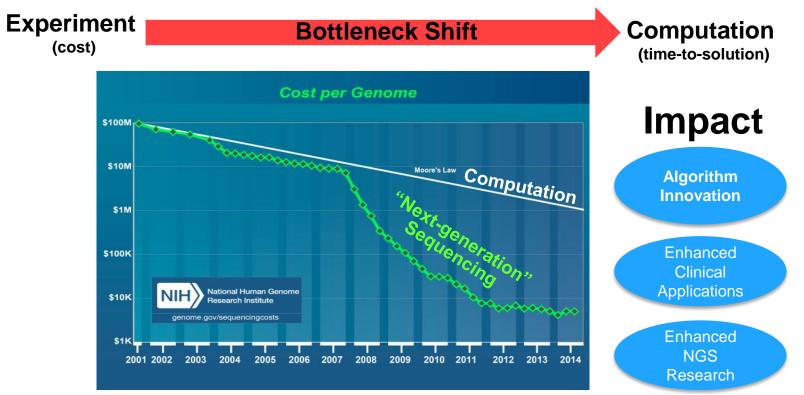Most applications are developed as a single executable

↔

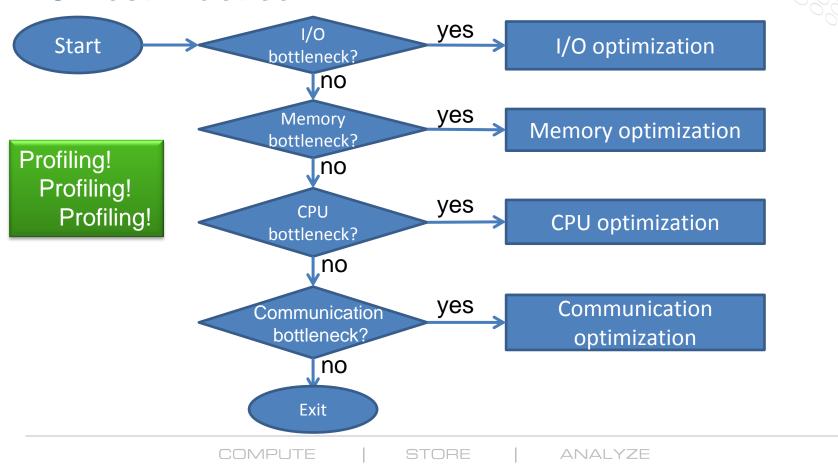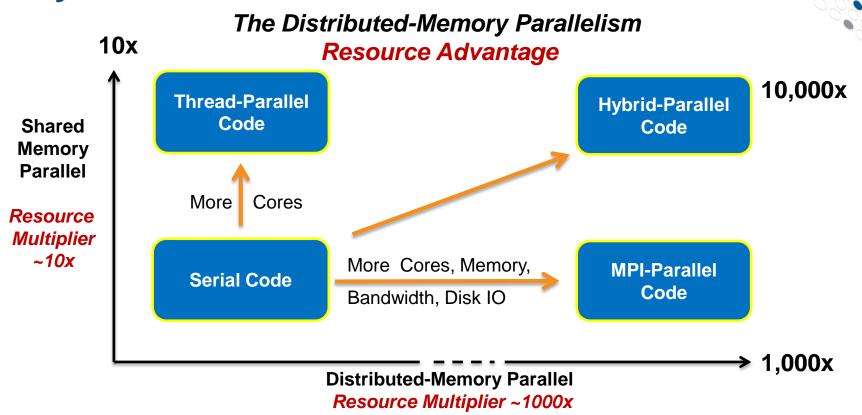NGS relies on running a collection of applications via a workflow

# NGS Bottlenecks are Now Computational

**Experiment**
(cost)

→ **Bottleneck Shift** →

**Computation**
(time-to-solution)



**Cost per Genome**

- $100M
- $10M — Moore's Law — Computation
- $1M
- $100K — "Next-generation Sequencing"
- $10K
- $1K

NIH — National Human Genome Research Institute — genome.gov/sequencingcosts

2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014

## Impact

- **Algorithm Innovation**
- Enhanced Clinical Applications
- Enhanced NGS Research

# HPC Best Practice

# Why HPC Best Practice?



**The Distributed-Memory Parallelism**
*Resource Advantage*

10x

**Shared Memory Parallel**

*Resource Multiplier ~10x*

| Thread-Parallel Code | | Hybrid-Parallel Code | 10,000x |

More Cores

Serial Code

More Cores, Memory, Bandwidth, Disk IO

MPI-Parallel Code

1,000x

**Distributed-Memory Parallel**
*Resource Multiplier ~1000x*

# TrinityRNA-Seq – How it works:



**RNA-Seq "reads"** → **Linear "contigs"** → **de-Bruijn graphs** → **Transcripts + Isoforms**

>a121:len=5845
>a122:len=2560
>a123:len=4443
>a124:len=48
>a125:len=8876
>a126:len=66

Thousands of <u>disjoint</u> graphs

...CTTCGCAA...TGATCGGAT...
...ATTCGCAA...TCATCGGAT...

Short reads

illumina®

**https://github.com/trinityrnaseq/trinityrnaseq/releases/tag/v2.0.6**

# TrinityRNA-Seq – How it works:



**Fastq file:**

**Fasta file:**

# MPI-Inchworm Algorithm

**Fastq file:**

**Fasta file:**

- ○ MPI_Send
- ○ MPI_Recv

Parallel read

**Phase 1: Distribute k-mers**

**Phase 2: Build contigs**

**Required few changes to the original code**

**Bowtie sort (linux)**

**GraphFromFasta**

**ReadsToTranscripts sort (linux)**

**QuantifyGraph**

**Butterfly**

COMPUTE | STORE | ANALYZE

# MPI-Inchworm Algorithm

Thread-Parallel Code

More Cores

Serial Code

More Cores, Memory, Bandwidth, Disk IO

Hybrid-Parallel Code

MPI-Parallel Code

- **Phase 1: Distribute k-mers**
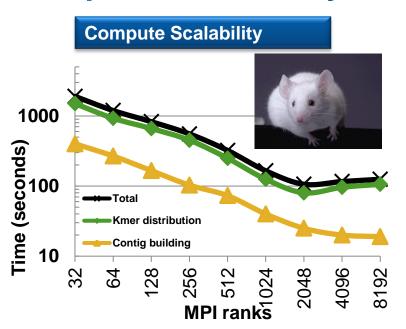
- **Phase 2: Build contigs**

# Hardware Resources

- Cray XC40, 64-bit Intel® Xeon® E5-2698 V3 "**Haswell**"

- **16 core** 2.3 GHz processor

- T**wo processors** per compute node and **384 processors per cabinet**

- The processor **peak performance** per core is 36.8 GF

- The memory consists of **128 GB** DDR4-**2133 MHz** per compute node

- **Memory bandwidth is 120 GB/s per node**

- The system interconnect is **Cray Aries multilevel dragonfly topology**

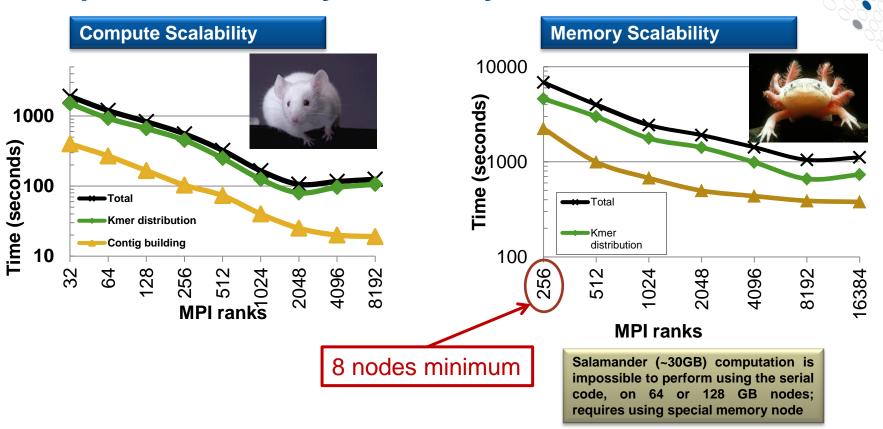> Essentially, 2 cabinets XC40 (2x192 nodes) is ideal for running MPI-inchworm (one run after another).

# Compute and Memory Scalability



Compute Scalability

Chart: Time (seconds) vs MPI ranks (32, 64, 128, 256, 512, 1024, 2048, 4096, 8192)
Legend: Total, Kmer distribution, Contig building

# Compute and Memory Scalability



8 nodes minimum

Salamander (~30GB) computation is impossible to perform using the serial code, on 64 or 128 GB nodes; requires using special memory node

# Compute and Memory Scalability



**Compute Scalability**

**Memory Scalability**

Distributed memory allows research on problems that otherwise would not be feasible

Salamander (~30GB) computation is impossible to perform using the serial code, on 64 or 128 GB nodes; requires using special memory node

# Conclusions

- **HPC Best Practice** can be applied to the **parallelization of NGS workflows.**
- The **distributed MPI-Inchworm** can now **utilize 4096 and more cores.**
- <u>Any</u> **bioinformatics** workflow can greatly **benefit from HPC.**
- Distributed-memory parallelism **eliminates** the need for **hybrid configurations** with large shared-memory nodes.

# Acknowledgements

- We thank **Cray Inc**. for the computing time on the **XC30/XC40 marketing machine.**

- We thank **Broad Institute principal investigator Aviv Regev** for generously supporting Trinity development efforts.

- We thank **Jessica Whited** at the **Brigham Regenerative Medicine Center** for access to the axolotl RNA-Seq data.

- Research reported in this publication was supported by the **National Cancer Institute** of the **National Institutes of Health (NIH) under Award Number 1U24CA180922-01.**

- The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

CODE: https://github.com/trinityrnaseq/trinityrnaseq/releases/tag/v2.0.6

CUG paper: http://biorxiv.org/content/early/2015/04/07/017665

**Pierre Carrier:** pcarrier@cray.com