# kallisto live demo

Páll Melsted

# Overview

- Downloading and installing kallisto

- Ingredients for a Coelacanth analysis

    - Nikaido *et al.* Genome Research, 2013

- Indexing

- Quantification

# Downloading and installing kallisto

- Download the source code and compile?

# Downloading and installing kallisto

Mac: use Homebrew science

- Install Homebrew from http://brew.sh

```
> brew tap homebrew/science
> brew install kallisto
```
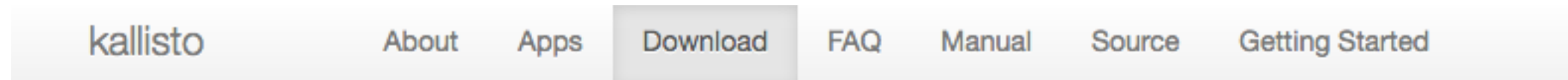
- Can also download binary, but Homebrew science adds a lot of useful tools

# Downloading and installing kallisto

Linux: download the binary (works on all distributions) from https://pachterlab.github.io/kallisto

# Downloading and installing kallisto

Windows: download binary and run using command line (cmd)

# Downloading and installing kallisto

Windows: download binary and run using command line (cmd)

# Requirements for running kallisto

- Reference transcriptome

- Paired end reads, FASTQ files

# Coelacanth analysis



- Reference transcriptome

- Paired end reads, FASTQ files

  - Downloaded from the short read archive (SRA)

    - 19 samples

      - 6 tissues - gill, kidney, pharynx, pectoral fin, pelvic fin, muscle

      - ~3 technical replicates per tissue

# Reference transcriptome

# Reference transcriptome

# Reference transcriptome

# Reference transcriptome



**Only download the transcriptome**

# Reference transcriptome



**Only download the transcriptome**

# Reference transcriptome

**Index of /pub/release-86/fasta/latimeria_chalumnae/**

| Name | Size | Date Modified |
|------|------|---------------|
| [parent directory] | | |
| cdna/ | | 9/26/16, 3:06:00 PM |
| cds/ | | 9/26/16, 3:06:00 PM |
| dna/ | | 9/26/16, 3:11:00 PM |
| ncrna/ | | 9/26/16, 3:11:00 PM |
| pep/ | | 9/26/16, 3:11:00 PM |

# Reference transcriptome

## Index of /pub/release-86/fasta/latimeria_chalumnae/

| Name | Size | Date Modified |
|------|------|---------------|
| [parent directory] | | |
| cdna/ | | 9/26/16, 3:06:00 PM |
| cds/ | | 9/26/16, 3:06:00 PM |
| dna/ | | 9/26/16, 3:11:00 PM |
| ncrna/ | | 9/26/16, 3:11:00 PM |
| pep/ | | 9/26/16, 3:11:00 PM |

## Index of /pub/release-86/fasta/latimeria_chalumnae/cdna/

| Name | Size | Date Modified |
|------|------|---------------|
| [parent directory] | | |
| CHECKSUMS | 134 B | 9/23/16, 2:47:00 PM |
| Latimeria_chalumnae.LatCha1.cdna.abinitio.fa.gz | 33.6 MB | 9/20/16, 3:05:00 PM |
| Latimeria_chalumnae.LatCha1.cdna.all.fa.gz | 16.1 MB | 9/20/16, 2:34:00 PM |
| README | 3.1 kB | 9/20/16, 3:05:00 PM |

# Indexing command

- Creates an index of the transcriptome reference that kallisto will use for quantification

    - Required only once per reference transcriptome

```
> time kallisto index -i index.idx
  Latimeria_chalumnae.LatCha1.cdna.all.fa.gz
```

# Indexing command

- Creates an index of the transcriptome reference that kallisto will use for quantification

  - Required only once per reference transcriptome

```
> time kallisto index -i index.idx
  Latimeria_chalumnae.LatCha1.cdna.all.fa.gz
```

where to store index

# Indexing command

- Creates an index of the transcriptome reference that kallisto will use for quantification

  - Required only once per reference transcriptome

```
> time kallisto index -i index.idx
  Latimeria_chalumnae.LatCha1.cdna.all.fa.gz
```

where to store index

Reference transcriptome

# Indexing guide

What value of k should be chosen?

- k-mer tradeoffs
  - high k - more specific
  - low k - robust to sequencing errors

# Indexing guide

What value of k should be chosen?

- k-mer tradeoffs

  - high k - more specific

  - low k - robust to sequencing errors

- general guide

  - for 75bp reads, use default, k=31

  - for 50bp reads, use default except if known issues

  - shorter reads, lower k=25 or k=21.

# Quantification

Quantification is run separately for each sample

```
> time kallisto quant -i index/index.idx
  -o results/DRR002318 -b 30 -t 4
  data/DRR002318_1.fastq.gz
  data/DRR002318_2.fastq.gz
```

# Quantification

Quantification is run separately for each sample

index constructed

```
> time kallisto quant -i index/index.idx
  -o results/DRR002318 -b 30 -t 4
  data/DRR002318_1.fastq.gz
  data/DRR002318_2.fastq.gz
```

# Quantification

Quantification is run separately for each sample

where to put results

index constructed

```
> time kallisto quant -i index/index.idx
-o results/DRR002318 -b 30 -t 4
data/DRR002318_1.fastq.gz
data/DRR002318_2.fastq.gz
```

# Quantification

Quantification is run separately for each sample

where to put results

index constructed

```
> time kallisto quant -i index/index.idx
  -o results/DRR002318 -b 30 -t 4
  data/DRR002318_1.fastq.gz
  data/DRR002318_2.fastq.gz
```

number of bootstraps

# Quantification

Quantification is run separately for each sample

where to put results

index constructed

```
> time kallisto quant -i index/index.idx
  -o results/DRR002318 -b 30 -t 4
  data/DRR002318_1.fastq.gz
  data/DRR002318_2.fastq.gz
```

number of bootstraps

threads used

# Quantification

Quantification is run separately for each sample
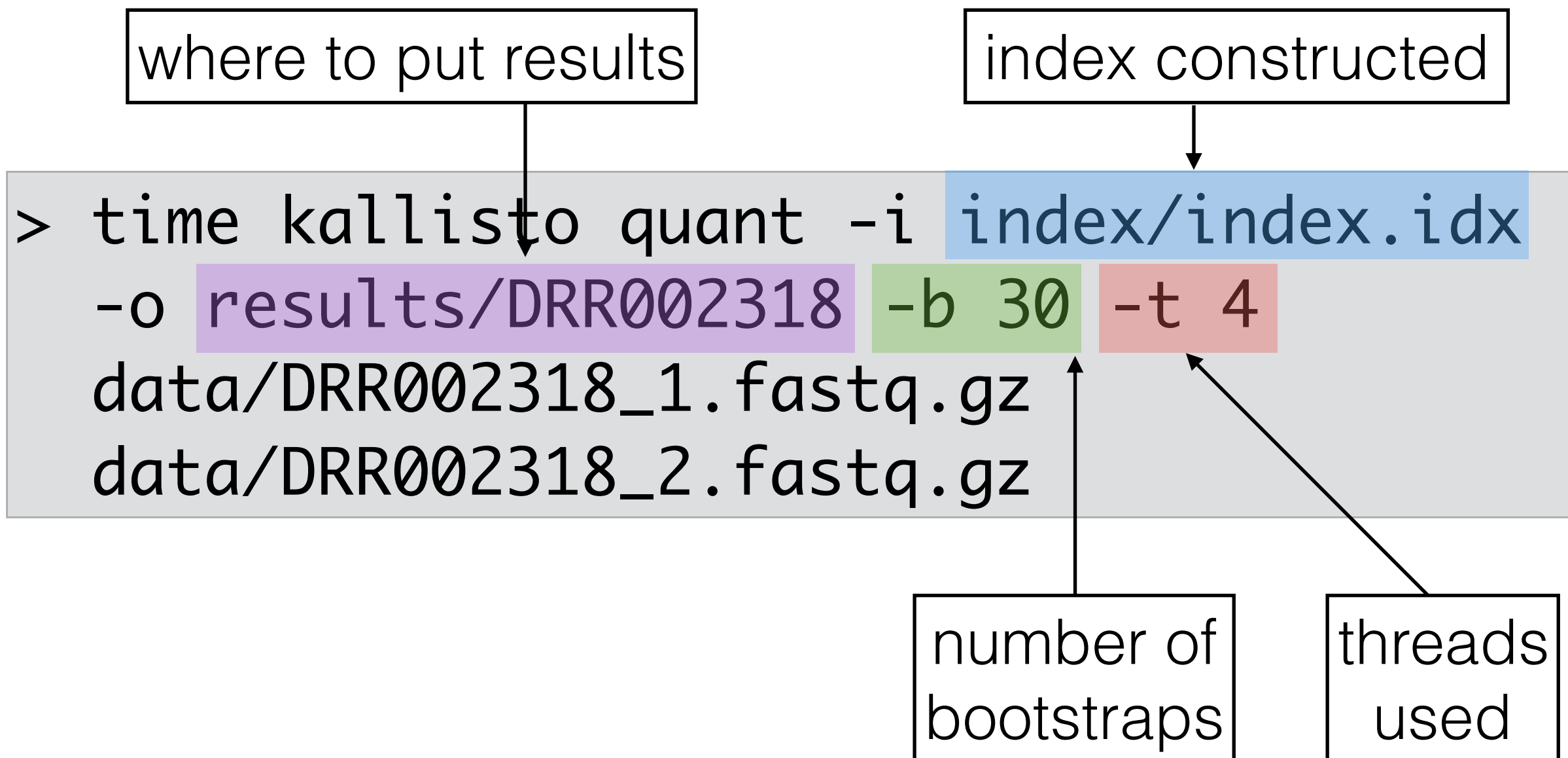
where to put results

index constructed

```
> time kallisto quant -i index/index.idx
 -o results/DRR002318 -b 30 -t 4
  data/DRR002318_1.fastq.gz
  data/DRR002318_2.fastq.gz
```
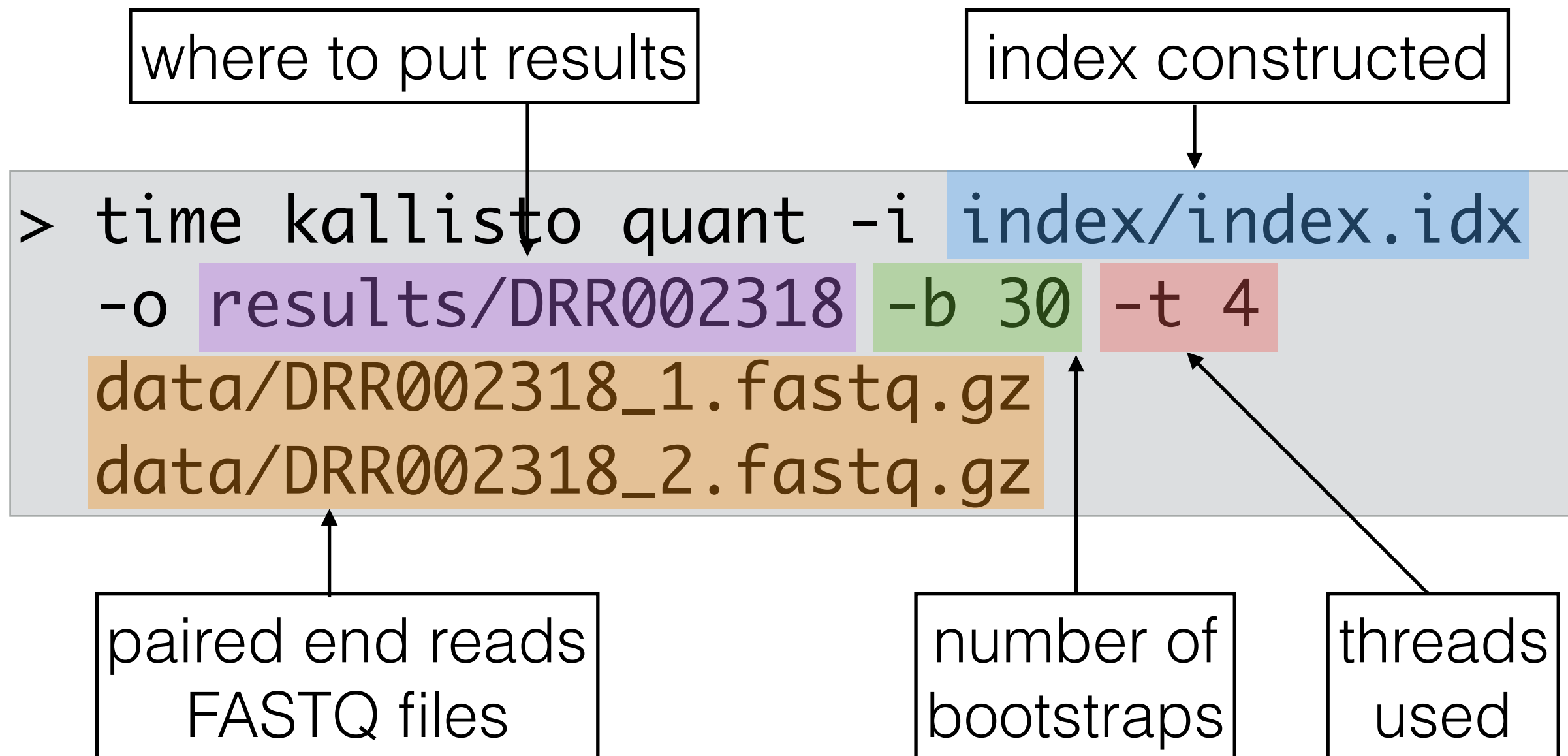
paired end reads
FASTQ files

number of
bootstraps

threads
used

# Results

- The output is stored in the specified directory

- `abundance.h5` - HDF5 compressed file, not human readable

  - contains the quantifications, all bootstraps and other information

- `abundance.tsv` - Tab Separated file

  - abundances and counts for each transcript

- `run_info.json` - JSON formatted file

  - information about the run