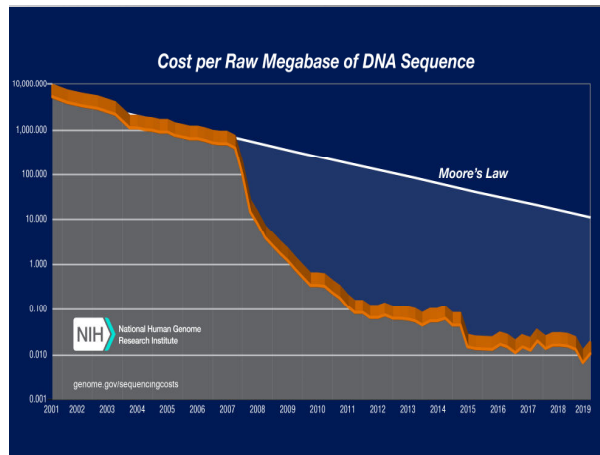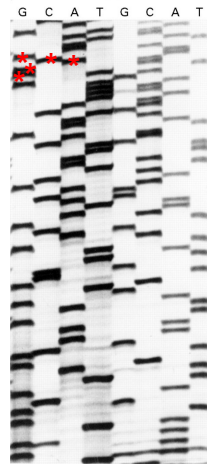1

## Next Generation Sequencing Applications

- The ability to generate large amounts of nucleotide sequence data has revolutionized biology in the past decade

- Applications include the *de novo* sequencing of genomes, transcriptomes, metagenomes, protein-genome interactions, etc.

- Advances driven by the promise of personalized genomic medicine, X Prize ($10 million reward to sequence 100 genomes in 30 days for $1000 each*), etc.

\* this X Prize was actually canceled because it was "outpaced by innovation" in sequencing technology, see https://en.wikipedia.org/wiki/Archon_X_Prize. The $1000 genome was achieved at the start of 2017…

🌿 **Bioinformatic Methods I**
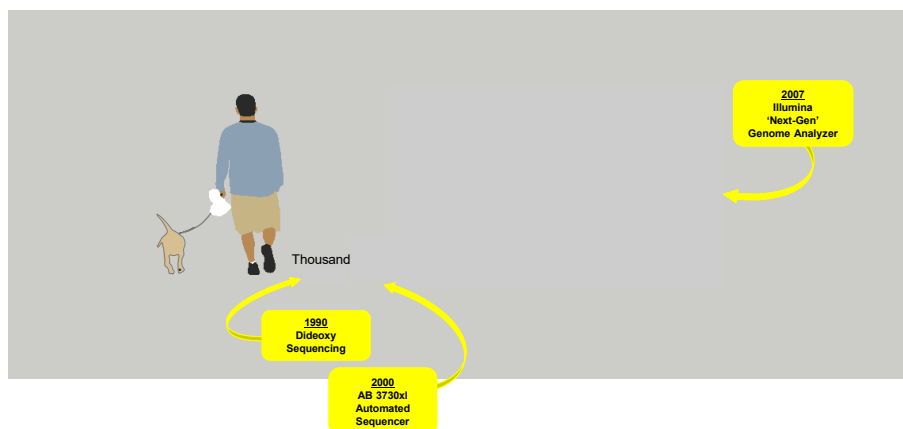
2

## Sequencing Rates and Costs

**Bioinformatic Methods I**

N. Provart & D. Guttman · Intro for Lab 6 · Slide 3
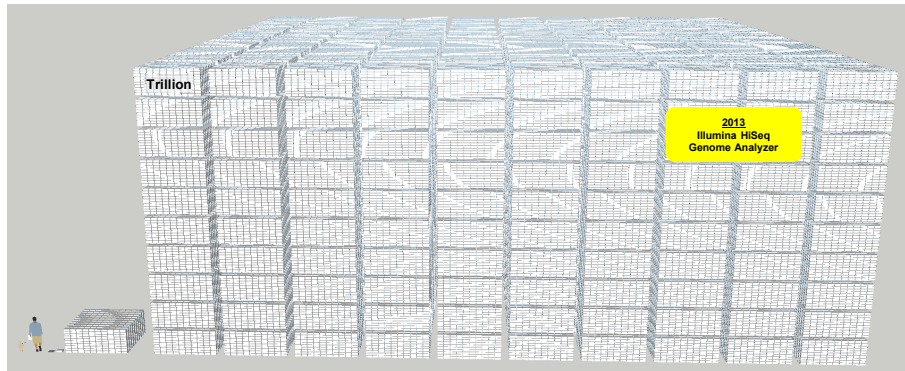
3

## Sequencing Rates and Costs



**Bioinformatic Methods I**

N. Provart & D. Guttman · Intro for Lab 6 · Slide 4
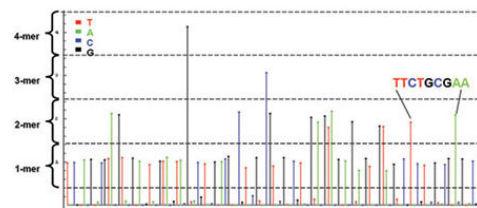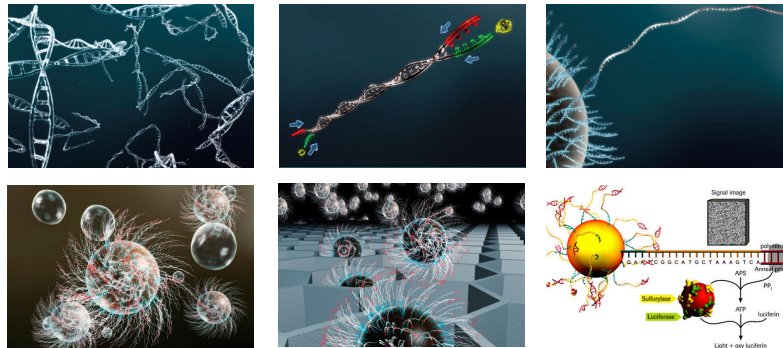
4

## Sequencing Rates and Costs



Trillion

**2013**
**Illumina HiSeq**
**Genome Analyzer**

🧪 **Bioinformatic Methods I**

5

## Sequencing Rates and Costs

| Human Gut Microbiome | | Sanger |
|---|---|---|
| Number of Species | 1000 | |
| Average Genome Size | 3 Mb | |
| Microbiome Size | 3 Gb | |
| Desired Coverage | 300 x | |
| Amount of Data Needed | ~ 1 Tb | |
| Read Length | | 750bp |
| Number of **Runs** Needed | | ~14M |
| Cost | | $667B |

| Platform | Reads | Read Length (bases) | Paired Ends | Run Time (days) | Yield (Gb) | Rate (days/Gb) |
|---|---|---|---|---|---|---|
| **Sanger** | 96 | 750 | No | 0.5 | 0.00007 | ~7000 days |

🧪 **Bioinformatic Methods I**

6

## Roche 454 Pyrosequencing



Images from 454.com website

**Bioinformatic Methods I**

## Illumina Sequencing by Synthesis



Images from Illumina.com website

**Bioinformatic Methods I**

**Illumina Sequencing by Synthesis**



Images from Illumina.com website

🎓 **Bioinformatic Methods I**

N. Provart & D. Guttman · Intro for Lab 6 · Slide 9

9

---

**Next Generation Sequencing Technologies**

- **Illumina GAllx, HiSeq2000 + newer machines (short reads, sequencing by synthesis)**
- **Roche 454 FLX, GS Junior (short reads, pyrosequencing)** – *now discontinued*
- **ABI SOLiD (short reads, sequencing by ligation)** – *good for SNP calling*

- PacBio (long reads*, real time single molecule sequencing)
- Oxford Nanopore (long reads*, nanopore sequencing)

- MGI DNBSEQ (short reads, DNA nanoball rolling circle sequencing)

> \* Long reads can really help with the assembly of complex eukaryotic genomes and transcriptomes, in a manner analogous to paired-end sequencing on slide 19, but in the case of these reads the intervening sequence is known!
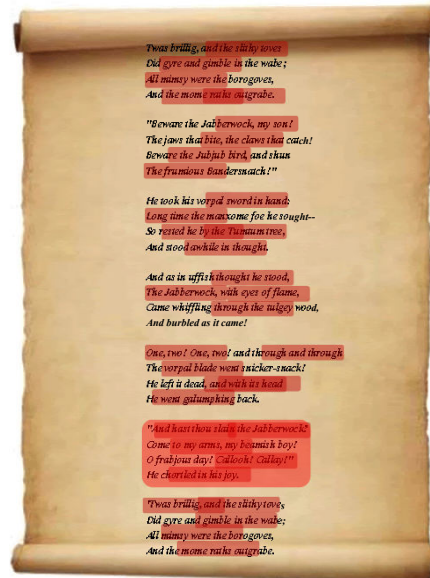
🎓 **Bioinformatic Methods I**

N. Provart & D. Guttman · Intro for Lab 6 · Slide 10

10

## NGS Assembly

- Cut poem into short "reads"

- Each part of poem is read multiple times

- Assemble poem by assembling overlapping reads

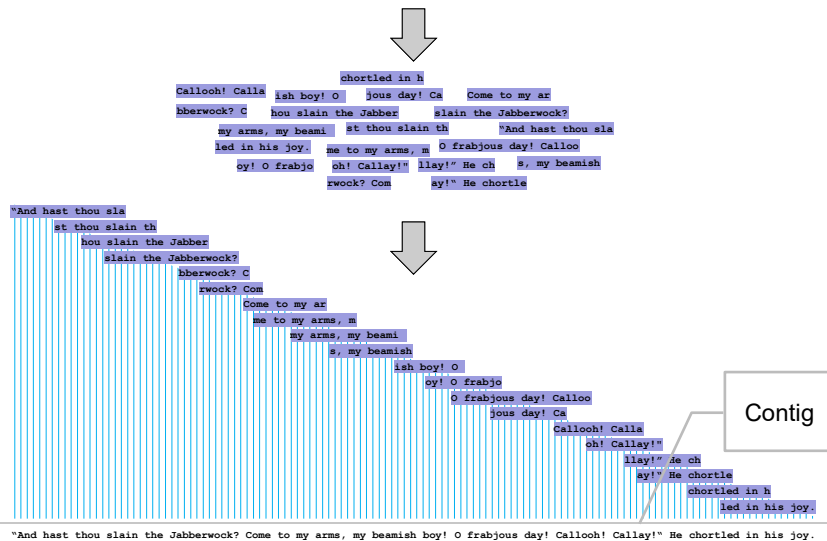"Jabberwocky" by Lewis Carroll, 1872

🌱 **Bioinformatic Methods I**

## Assembly

"And hast thou slain the Jabberwock? Come to my arms, my beamish boy! O frabjous day! Callooh! Callay!" He chortled in his joy.

Contig

"And hast thou slain the Jabberwock? Come to my arms, my beamish boy! O frabjous day! Callooh! Callay!" He chortled in his joy.
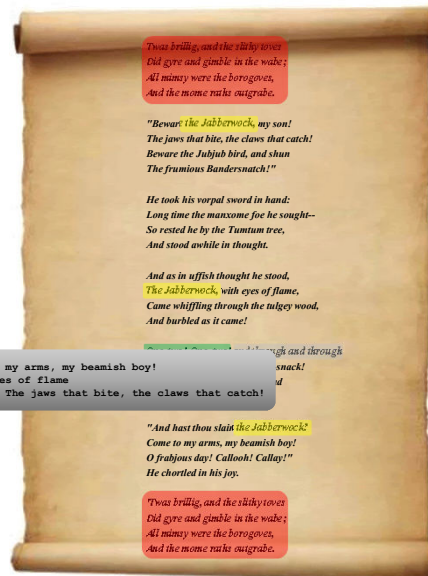
🌱 **Bioinformatic Methods I**

## Assemby – Strategy

- Cut poem into short "reads"

- Each part of poem is read multiple times

- Assemble poem by assembling overlapping reads

Problem
- Repetitive regions



```
          "And hast thou slain          ? Come to my arms, my beamish boy!
And as in uffish thought he stood,   the Jabberwock  , with eyes of flame
          "Beware                    , my son! The jaws that bite, the claws that catch!
```

## Assembly – Challenges

- 50 - 1000 Epic Poems (for human microbiome)

- Each poem is 500,000 – 10,000,000 characters long

- Each poem is written in an alphabet of only four letters

- Some poems are present in billions of copies, while others only present in only tens of copies

- Each poem is randomly cut up into "reads" ranging in length from 75 – 500 characters long

- There are as many as a billion such reads present in no particular order
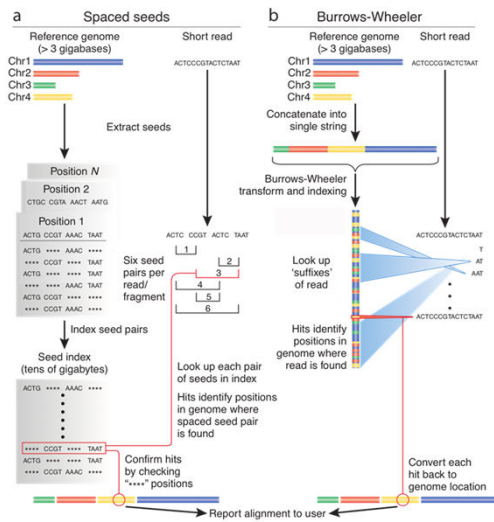
## Mapping Reads to Genomes - Overview



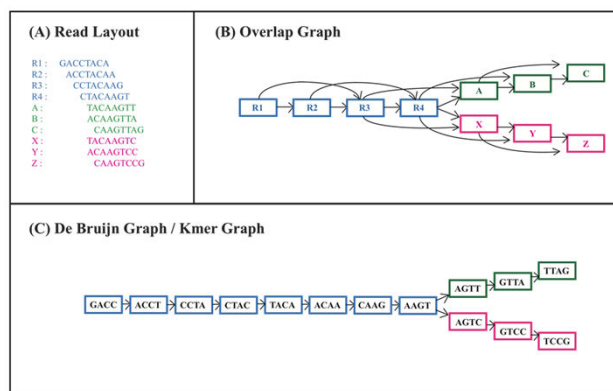Image from Trapnell & Salzberg (2009) Nature Biotech. 27:455-457.

## *De novo* Genome Assembly



Image courtesy of Jessica Yang, M.Sc. Thesis (2011), adapted from Schatz et al. (2010) Genome Res. 20:1165-73.

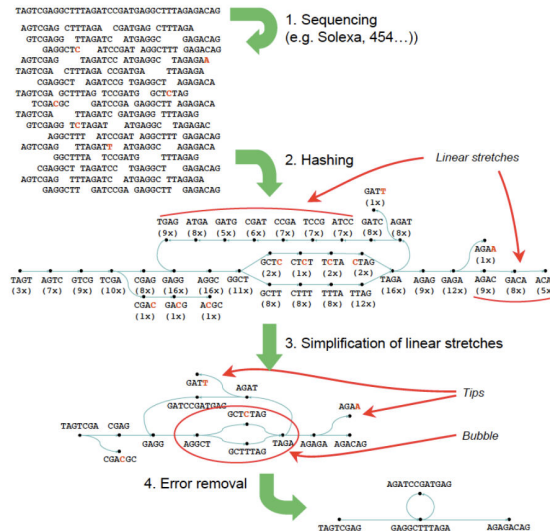### *De novo* Genome Assembly, e.g. Velvet



Image courtesy of Jim Noonan, Yale University.

**Bioinformatic Methods I**

N. Provart & D. Guttman · Intro for Lab 6 · Slide 7

17

---

### Assembly Quality Assessment Scores

- Typically, several assemblers (such as ABySS, SOAPdenovo, Velvet etc.) are used and then a final assessment of quality is made using parameters such as:

- Number of Contigs…generally, the fewer the better

- N50: the maximum length *L* such that 50% of all bases lie in contigs at least *L* bases long…generally, the longer the better

- Coverage: the number of reads covering each base in a contig

Note that *de novo* genome (and transcriptome) assembly, especially from short reads, is a very active area of research. The first "assemblathon" was recently held to assess many different methods.

Assemblathon 1: Earl et al. (2011) Genome Research. DOI: 10.1101/gr.126599.111.
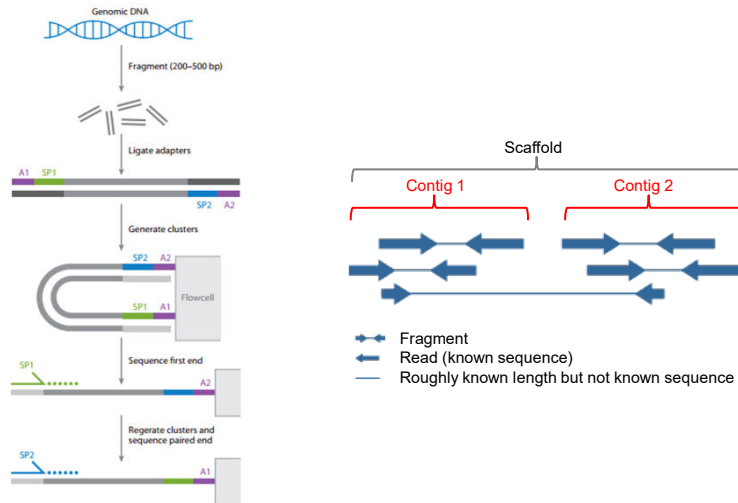
**Bioinformatic Methods I**

N. Provart & D. Guttman · Intro for Lab 6 · Slide 8

18

9

## How Paired-End Reads Can Facilitate Assembly
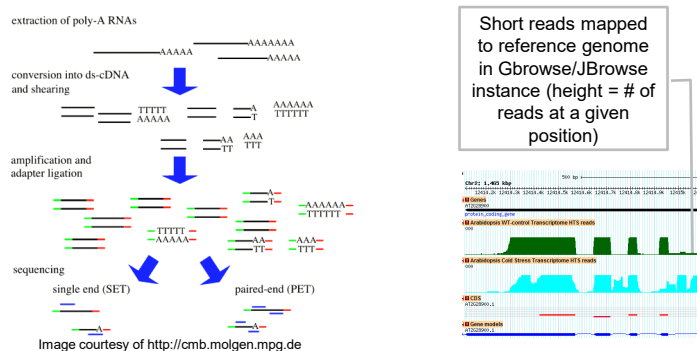


From Illumina.com website and http://genome.jgi.doe.gov/help/scaffolds.html

🛡 **Bioinformatic Methods I**

## RNA-Seq

RNA-seq is a powerful method for quantifying steady-state mRNA expression levels and detecting alternative splicing events in transcriptomes.

An RNA-seq pipeline involves creating cDNA, shearing, adding adapters, NGS, mapping reads, and summarizing the read counts (e.g. FPKM – fragments per kilobase per million fragments mapped). Evaluation of alternative splicing events may be visualized in a genome browser.
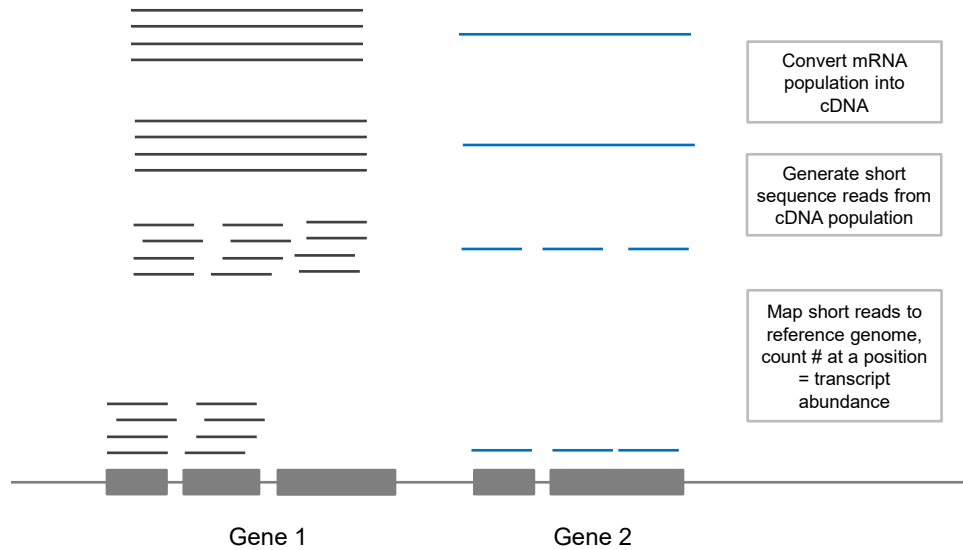


Image courtesy of http://cmb.molgen.mpg.de

Short reads mapped to reference genome in Gbrowse/JBrowse instance (height = # of reads at a given position)

🛡 **Bioinformatic Methods I**

## RNA-Seq can be used to identify alternative splicing events



Convert mRNA population into cDNA

Generate short sequence reads from cDNA population

Map short reads to reference genome, count # at a position = transcript abundance

Gene 1          Gene 2

🌱 **Bioinformatic Methods I**

21

## Alternative splicing is a common occurrence

"Of the 25,800 genes expressed along the [maize] leaf gradient we detected evidence of alternative splicing at 9,492 genes…Only 20,999 of [these] contain introns, so 56.4% of all possible targets showed evidence of alternative splicing."
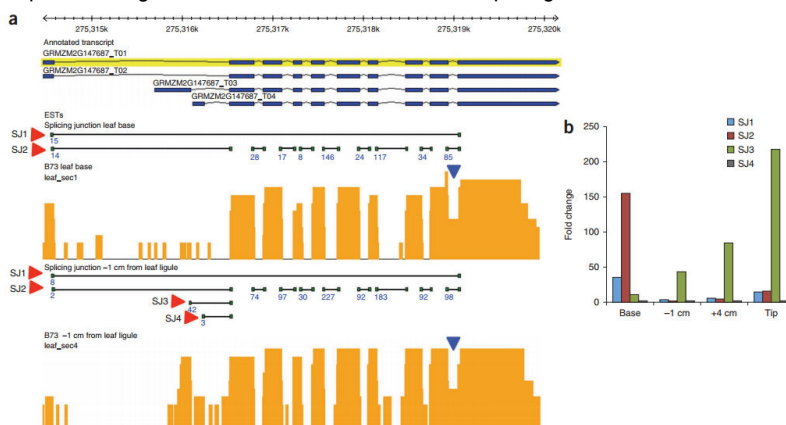


**Figure 3** Alternative splicing of GRMZM2G147687. (**a**) Genome browser shows the alignment of reads to splice junctions (green) and exons (yellow) of GRMZM2G147687. Red arrows indicate the alternative splice junctions within the first exon, and blue arrows indicate putative intron retention events (see also **Supplementary Fig. 3**). (**b**) Results of qRT-PCR showing accumulation of four isoforms along the developmental gradient.

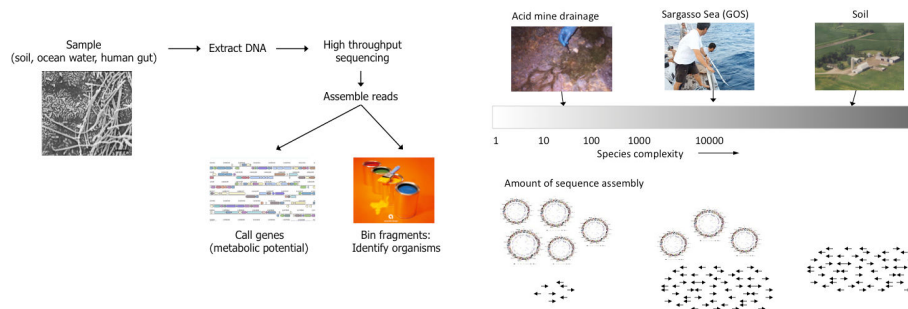Li~Brutnell (2010). Nature Genetics 42: 1060-1067

🌱 **Bioinformatic Methods I**

22

## Metagenomics

"The application of modern genomics techniques to the study of communities of microbial organisms directly in their natural environments, *bypassing the need for isolation and lab cultivation of individual species*"



Definition: Chen & Pacter (2005) PLoS CompBio 1:e24. Images courtesy of JGI.

23

---

## Some amazing things are possible with NGS!



ARTICLE

https://doi.org/10.1038/s41467-019-13549-9    OPEN

# A 5700 year-old human genome and oral microbiome from chewed birch pitch

Theis Z.T. Jensen [1,2,10], Jonas Niemann[1,2,10], Katrine Højholt Iversen [3,4,10], Anna K. Fotakis [1], Shyam Gopalakrishnan [1], Åshild J. Vågene[1], Mikkel Winther Pedersen [1], Mikkel-Holger S. Sinding [1], Martin R. Ellegaard [1], Morten E. Allentoft[1], Liam T. Lanigan[1], Alberto J. Taurozzi[1], Sofie Holtsmark Nielsen[1], Michael W. Dee[5], Martin N. Mortensen [6], Mads C. Christensen[6], Søren A. Sørensen[7], Matthew J. Collins[1,8], M. Thomas P. Gilbert [1,9], Martin Sikora [1], Simon Rasmussen [4] & Hannes Schroeder [1*]

24