# Genomic Data Compression and Processing for

# Large and Growing Databases
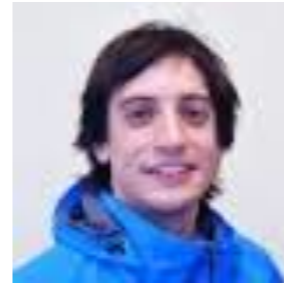
## Science of Information

**EE 25N**

**Tsachy Weissman**

# thanks

Idoia
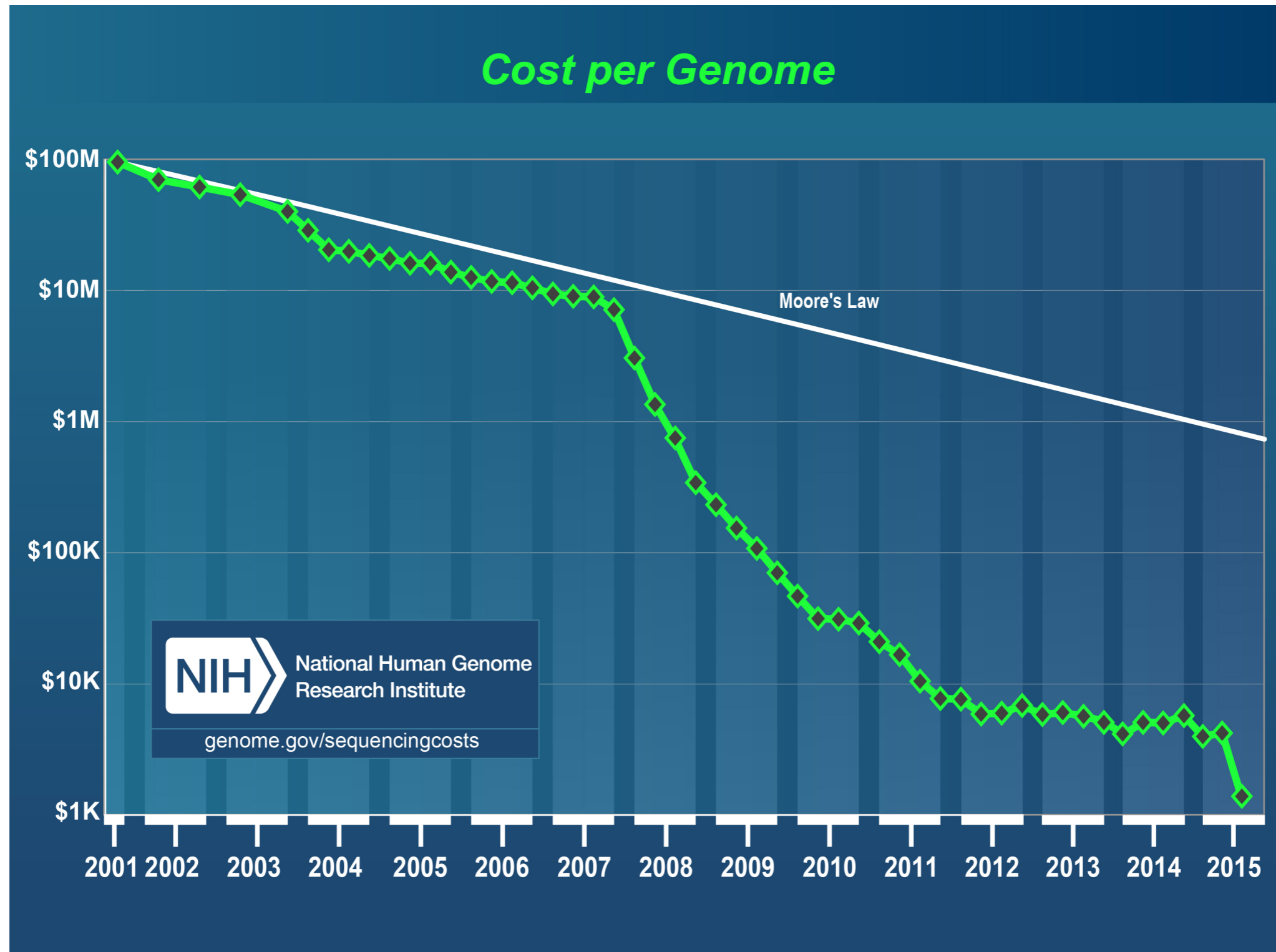Ochoa

Mikel
Hernaez

Shubham
Chandak

Kedar
Tatwawadi

Dmitri
Pavlichin

# a slide we've seen (multiple times..)



Cost per Genome

genome.gov/sequencingcosts

# why compression of genomic data?

- proxy for reduced cost of storage, communication, computation, processing, etc.
- compression as modelling

# 2 modes

- lossless
- lossy

(true) story of human genome compression

# Single Genome Compression

- **FASTA Compression**: Compression of a single genome

- Human genome can be represented using 2 bits/bp, compressed size $\approx 1GB$

- Specialized compressors: MFCompress[1] : $\approx 1.6$ bits/bp

- $H$(human genome): "Entropy" of the human genome

$$H(\text{human genome}) \lesssim 2 \text{ bits/bp} \sim 1 \text{ GB}$$

.

---

[1]A. J. Pinho and D. Pratas, "MFCompress: a compression tool for FASTA and multi-FASTA data", *Bioinformatics*, vol. 30, no. 1, pp. 117–118, 2014.

# Using a reference

- $H(\text{human genome}) \lesssim 1$ GB

- We can do better if we know another genome (reference)

- **Using a reference:** GenomeZip[2] compresses James Watson's genome using: $\approx 2.5MB$

$$H \left( \begin{array}{c} human \\ genome \end{array} \middle| \begin{array}{c} another \\ human \\ genome \end{array} \right) \lesssim 2.5MB$$

[2]D. S. Pavlichin, T. Weissman, and G. Yona, "The human genome contracts again", *Bioinformatics*, vol. 29, no. 17, pp. 2199–2202, 2013.

# Using a collection

- $H(\text{human genome}) \lesssim 1 \text{ GB}$

- $H\left(\begin{array}{l}\textit{human} \\ \textit{genome}\end{array}\middle|\begin{array}{l}\textit{another} \\ \textit{human} \\ \textit{genome}\end{array}\right) \lesssim 2.5 MB$

- **GTRAC**[3] compressor:

$$H\left(\begin{array}{l}\textit{human} \\ \textit{genome}\end{array}\middle|\begin{array}{c}1K \\ \textit{other} \\ \textit{genomes}\end{array}\right) \lesssim 1MB$$

- **GTC**[4] compressor:

$$H\left(\begin{array}{l}\textit{human} \\ \textit{genome}\end{array}\middle|\begin{array}{c}27K \\ \textit{other} \\ \textit{genomes}\end{array}\right) \lesssim 200KB$$

---

[3] K. Tatwawadi, M. Hernaez, I. Ochoa, *et al.*, "GTRAC: fast retrieval from compressed collections of genomic variants", *Bioinformatics*, vol. 32, no. 17, pp. i479–i486, 2016.

[4] A. Danek and S. Deorowicz, "GTC: how to maintain huge genotype collections in a compressed form", *Bioinformatics*, 2018.
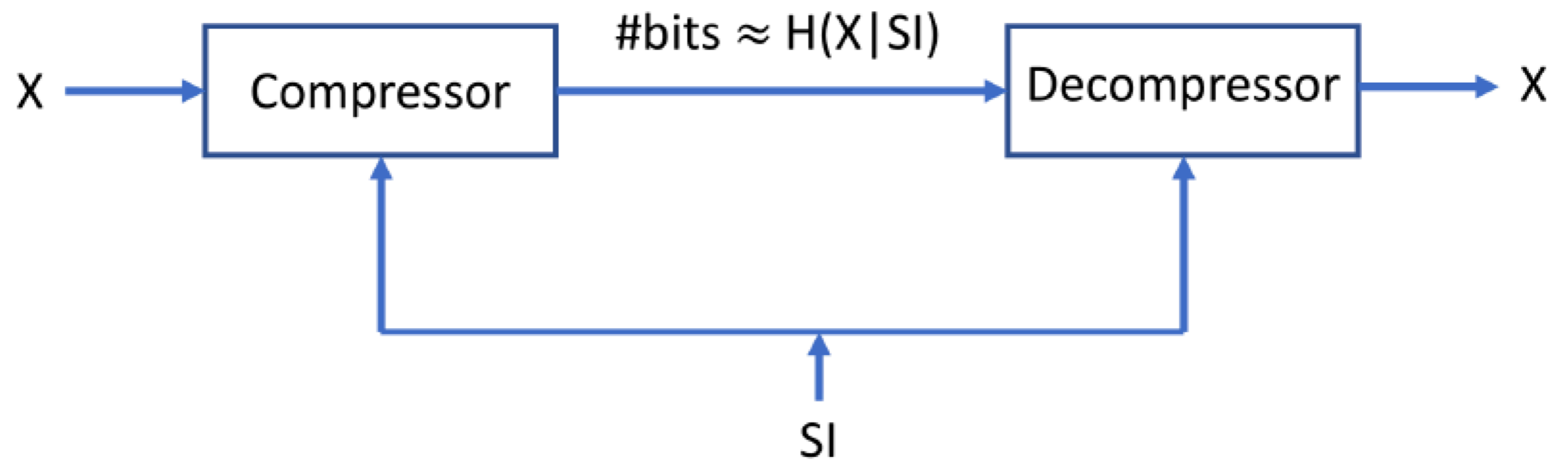
# Full Genome Compression

- $H(\text{human genome}) \lesssim 1 \text{ GB}$

- $H\left(\begin{smallmatrix} human \\ genome \end{smallmatrix} \middle| \begin{smallmatrix} another \\ human \\ genome \end{smallmatrix}\right) \lesssim 2.5MB$

- $H\left(\begin{smallmatrix} human \\ genome \end{smallmatrix} \middle| \begin{smallmatrix} 1K \\ other \\ genomes \end{smallmatrix}\right) \lesssim 1MB$

- $H\left(\begin{smallmatrix} human \\ genome \end{smallmatrix} \middle| \begin{smallmatrix} 27K \\ other \\ genomes \end{smallmatrix}\right) \lesssim 200KB$

# File Size per Genome Vs Database Size

# Information Theoretic Perspective

# Information Theoretic Perspective

# Information Theoretic Perspective



$$\text{X} \longrightarrow \boxed{\text{Compressor}} \xrightarrow{\#\text{bits} = H(X|SI)} \boxed{\text{Decompressor}} \longrightarrow \text{X}$$

SI

Slepian, David S.; Wolf, Jack K. (July 1973).
"Noiseless coding of correlated information sources"
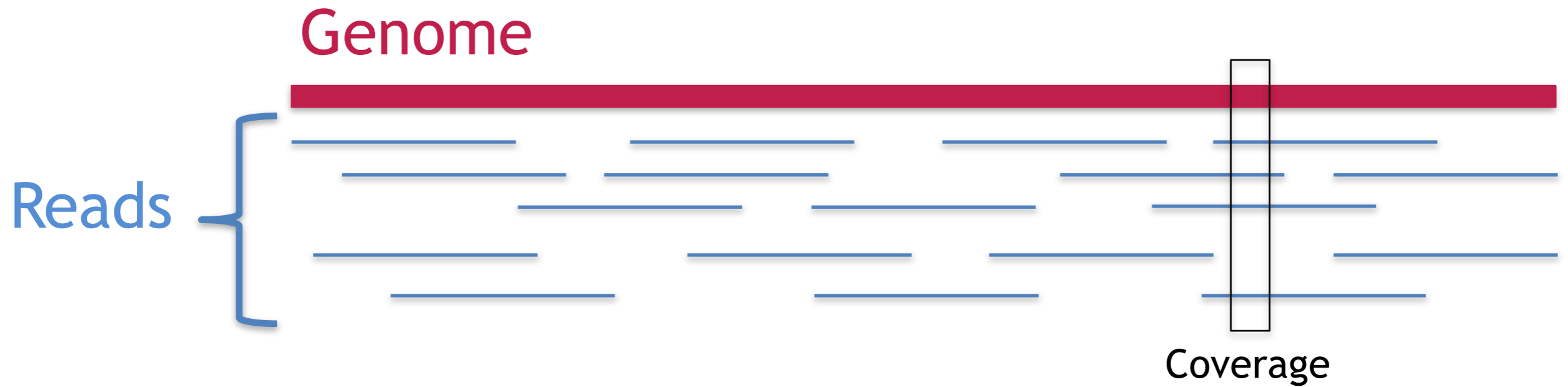
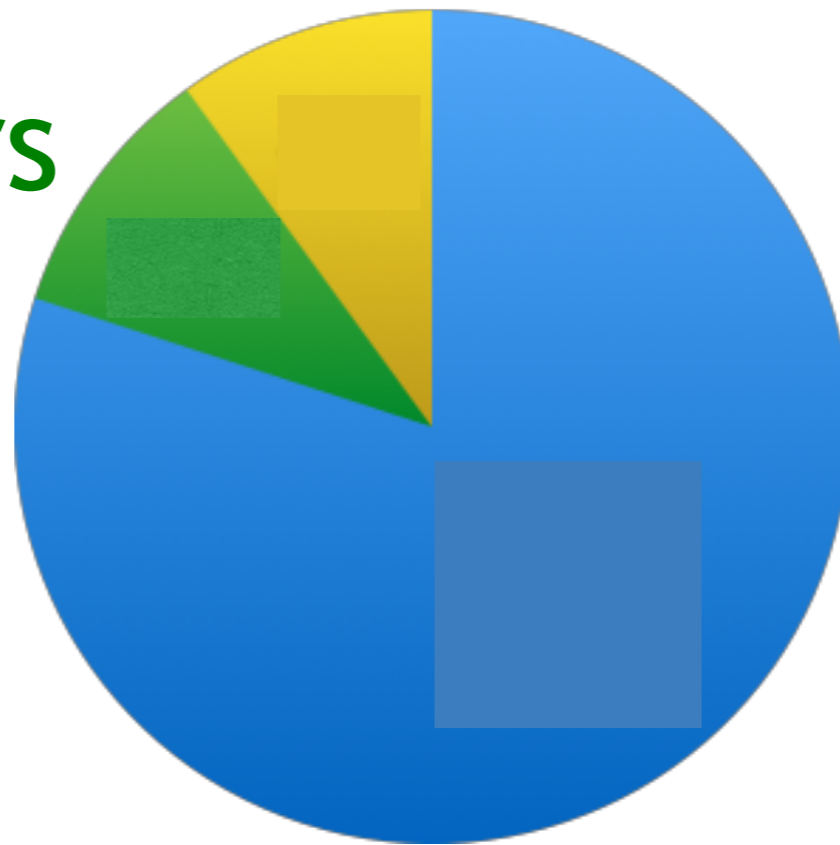# why lossy compression?
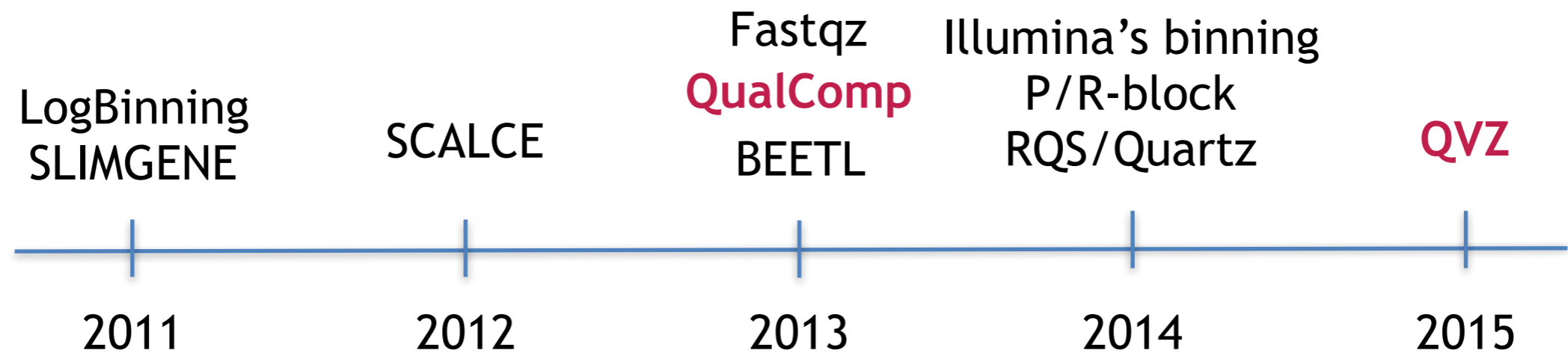
# why *lossy* compression of *genomic* data?

# genome sequencing



Genome

Reads

Coverage

Reads + alignment information

Identifiers

Quality scores

18

# lossy compressors of quality scores



**"QualComp: a new lossy compressor for quality scores based on rate distortion theory"**

**"QVZ: lossy compression of quality values"**

# how does
# lossy compression
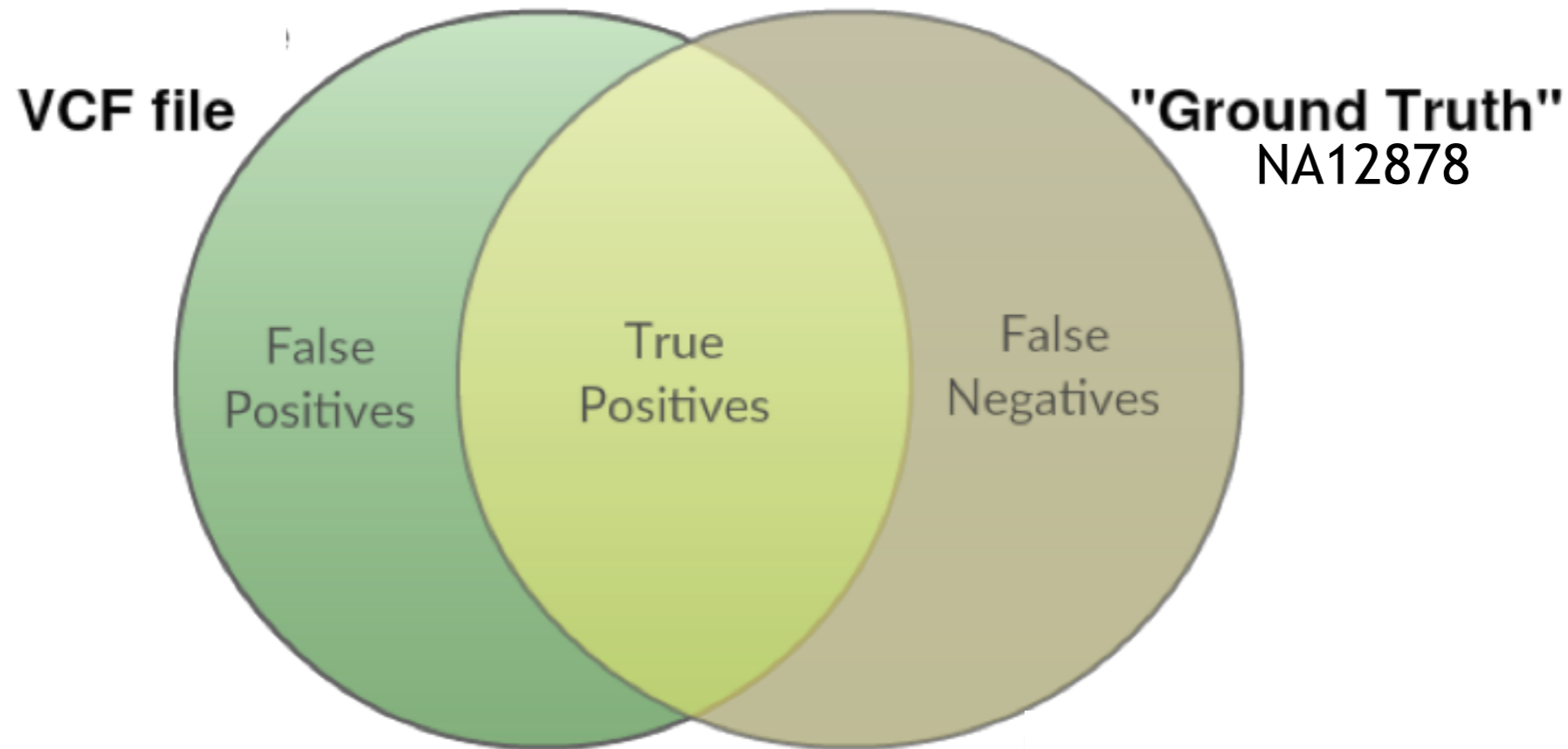# of quality scores
# affect the inference?

**"Effect of lossy compression of quality scores on variant calling"**

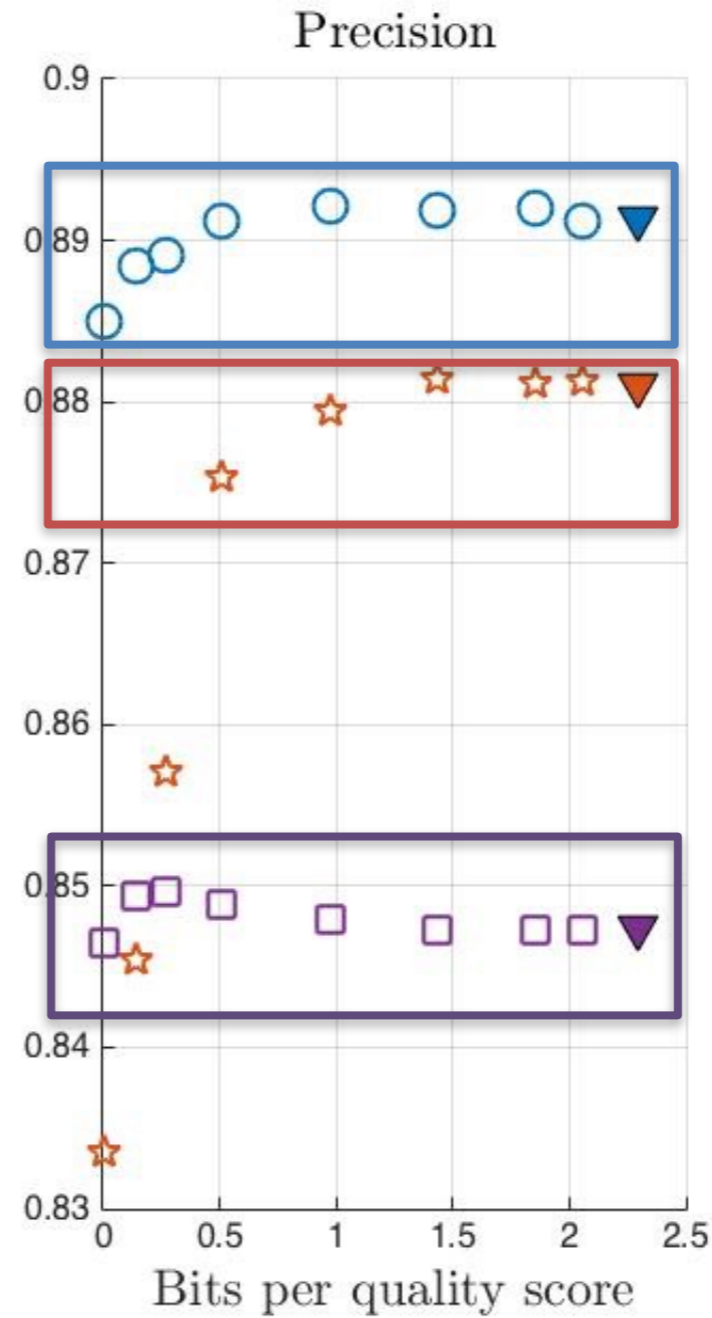with Idoia Ochoa, Mikel Hernaez, Rachel Goldfeder and Euan Ashley
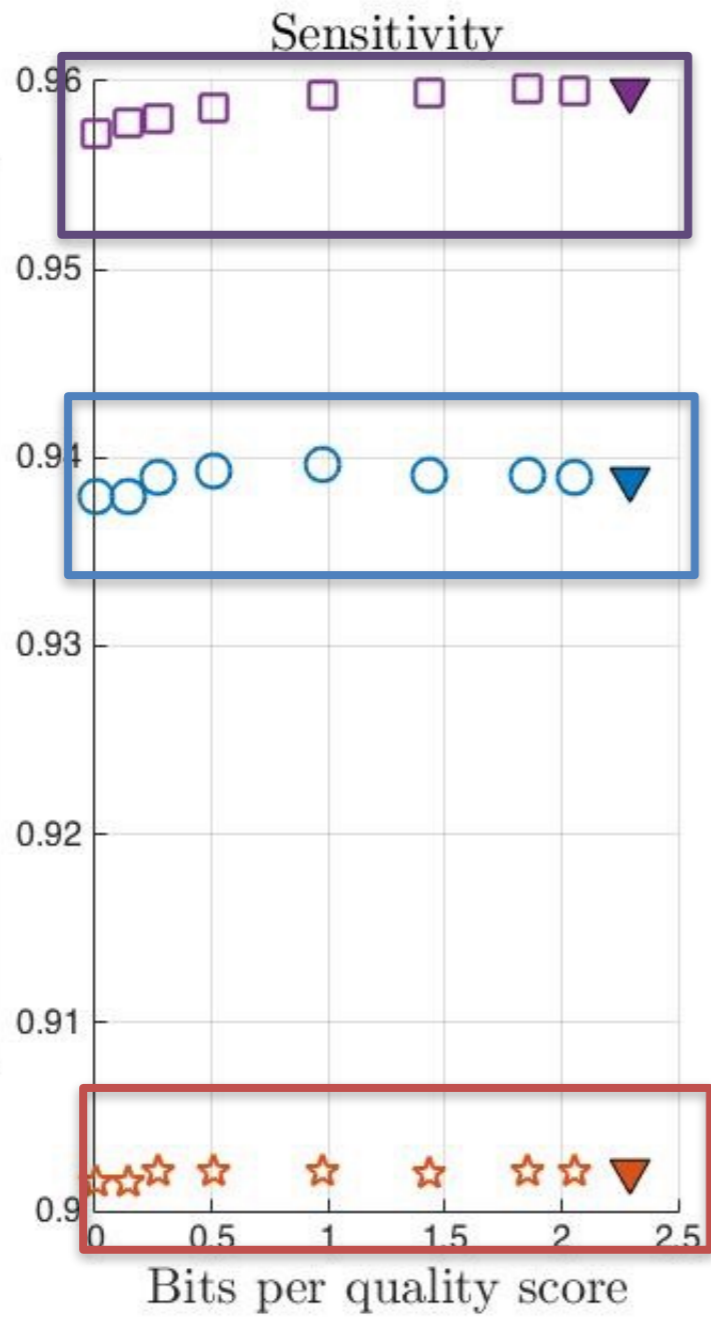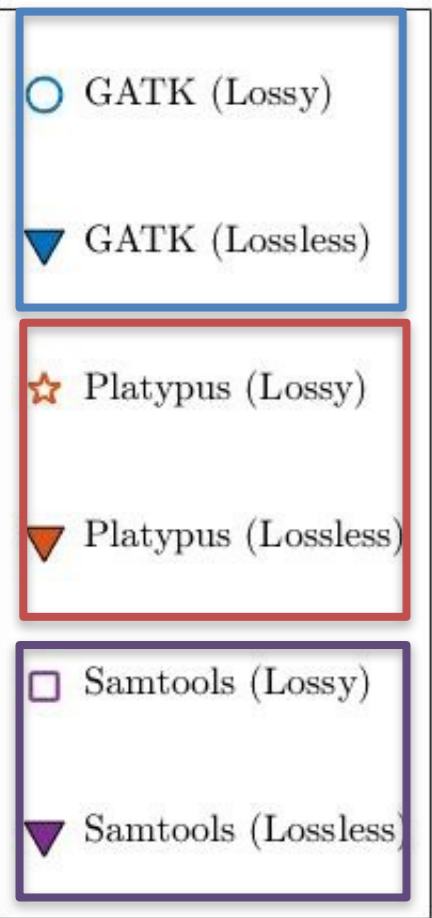
*Briefings in Bioinformatics, 2017*

Reference genome

Reference genome

FASTQ file → Alignment → SAM file → Variant caller → VCF file

- BWA
doesn't use the
Quality Scores
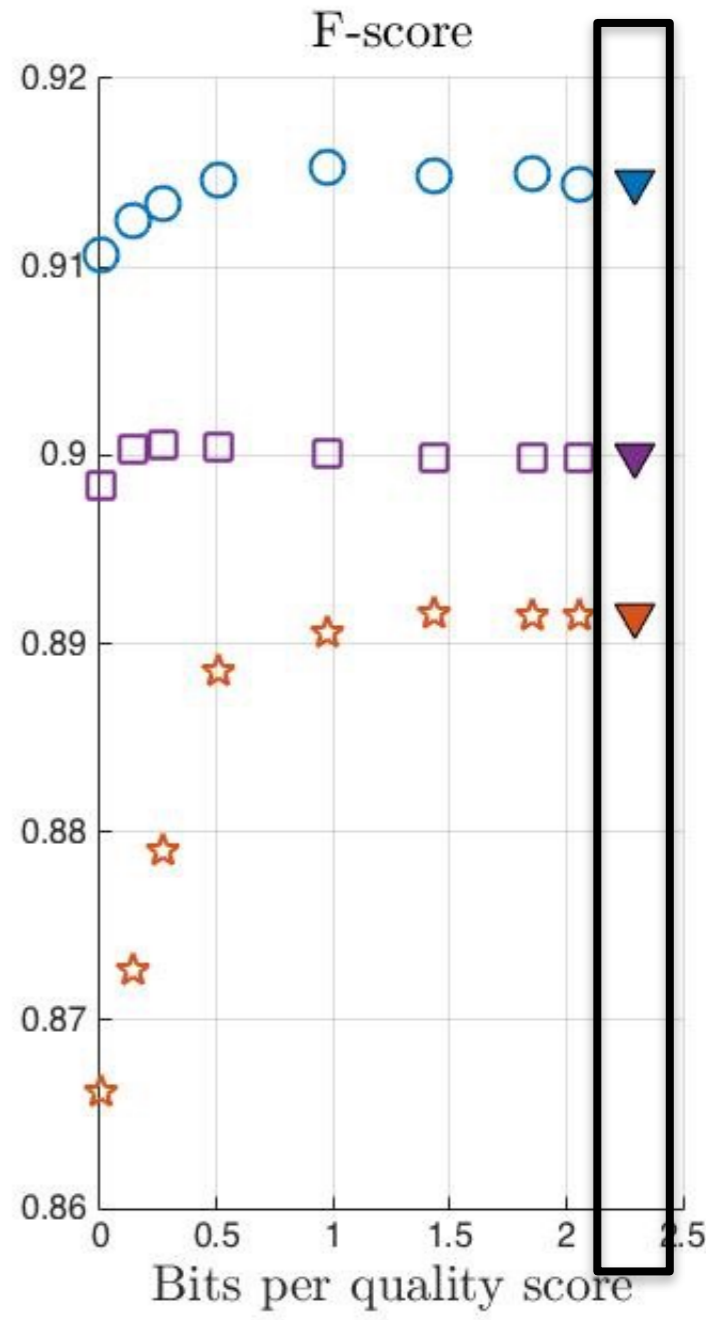
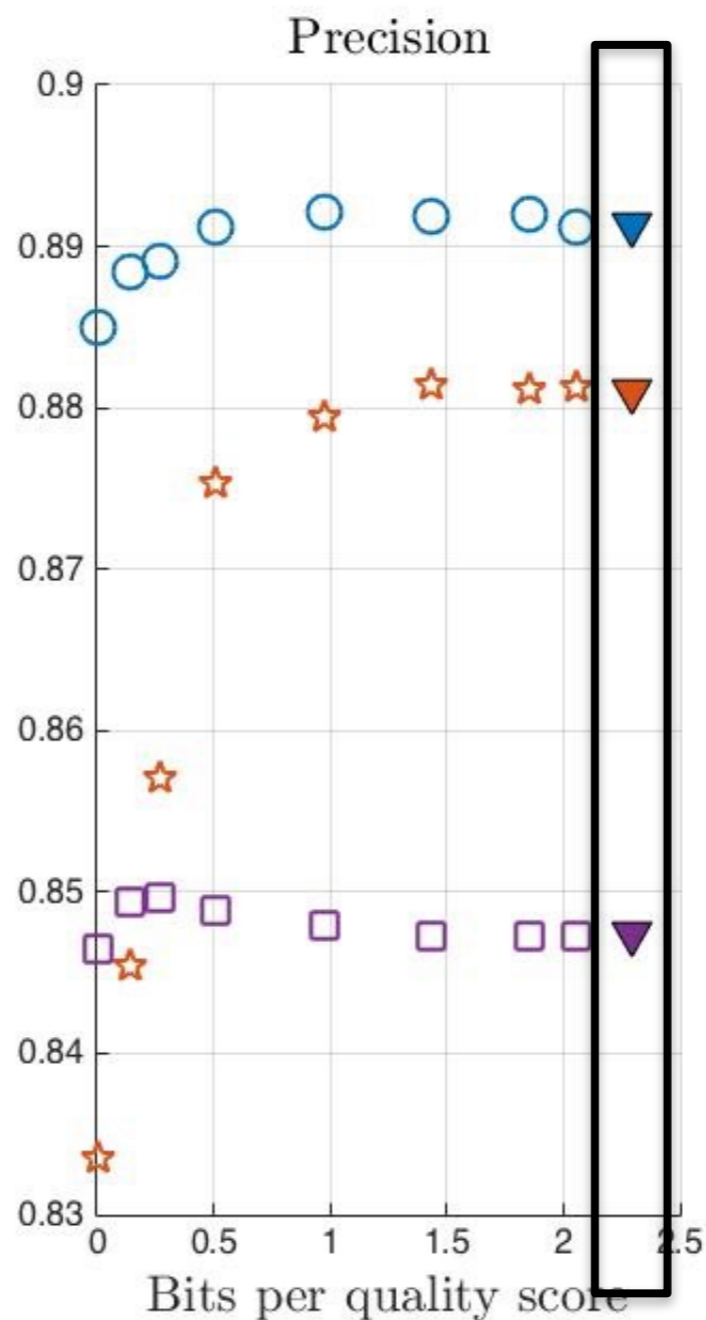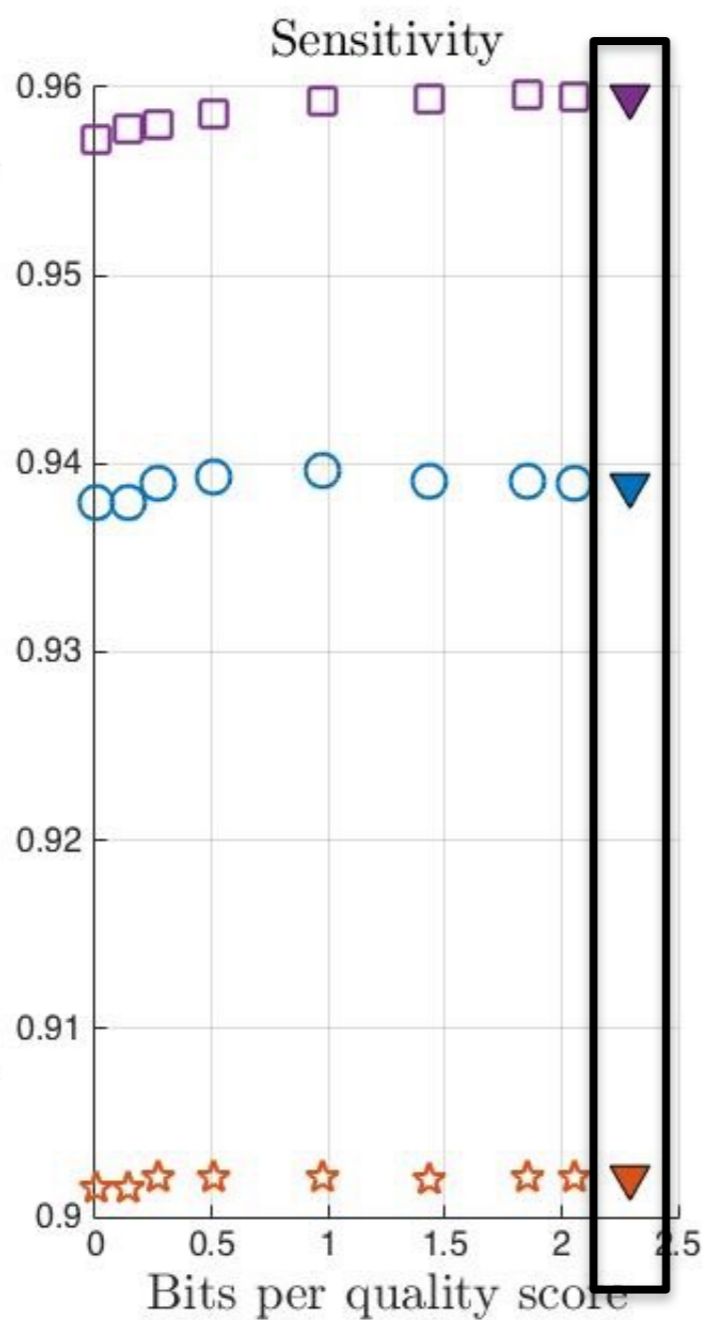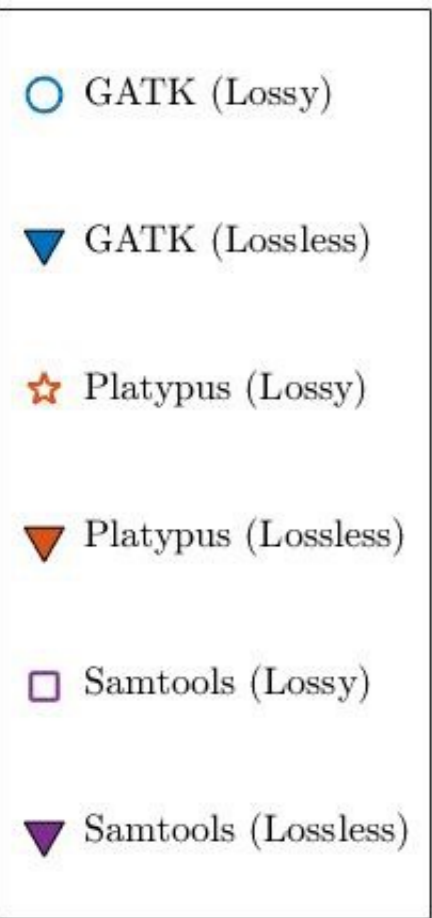→ Extract Quality Scores → Lossy Compression → Lossy Decompression →

# performance metrics



- *sensitivity*: T.P. / (T.P. + F.N.)
- p*recision*: T.P. / (T.P. + F.P.)
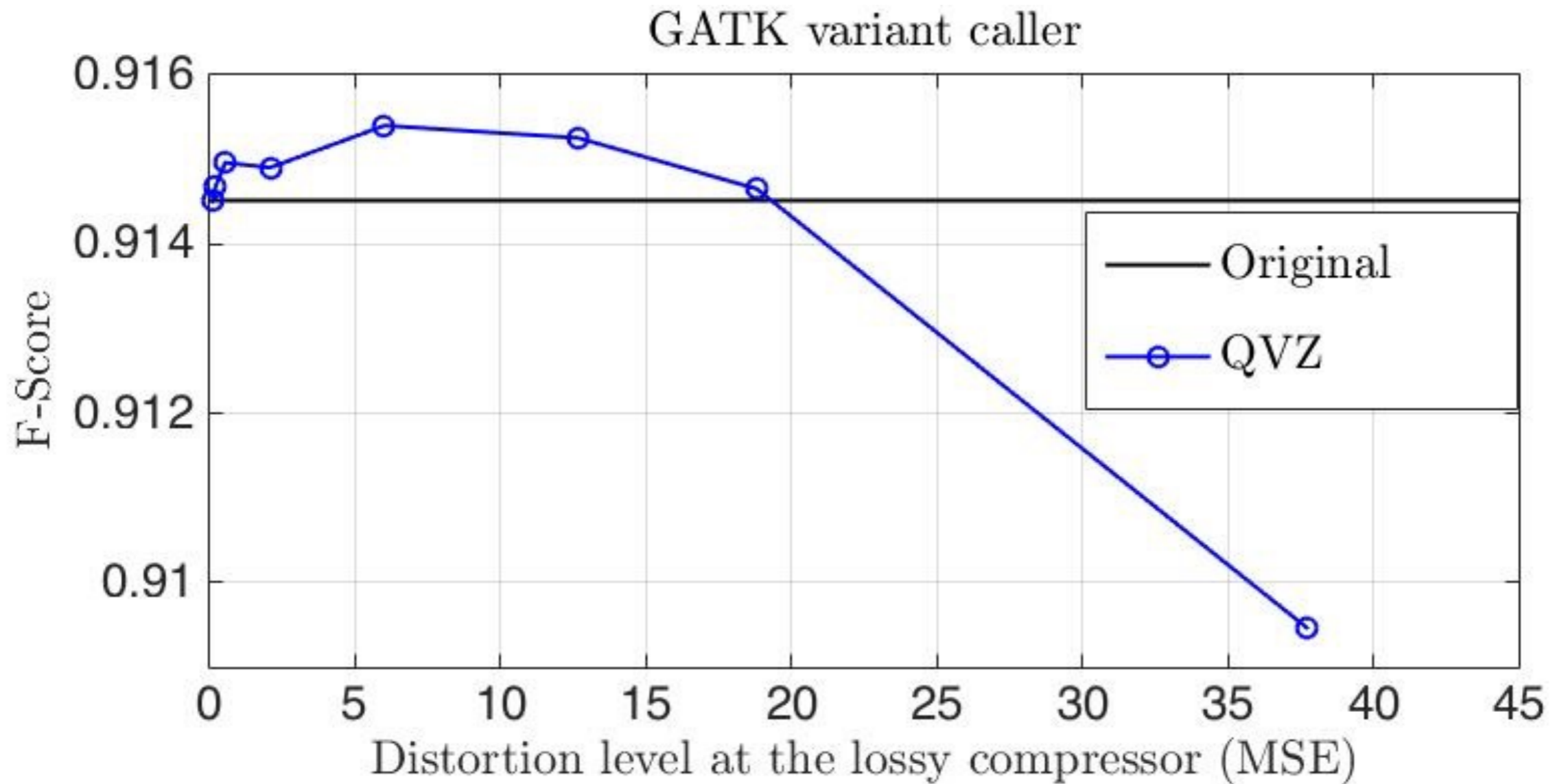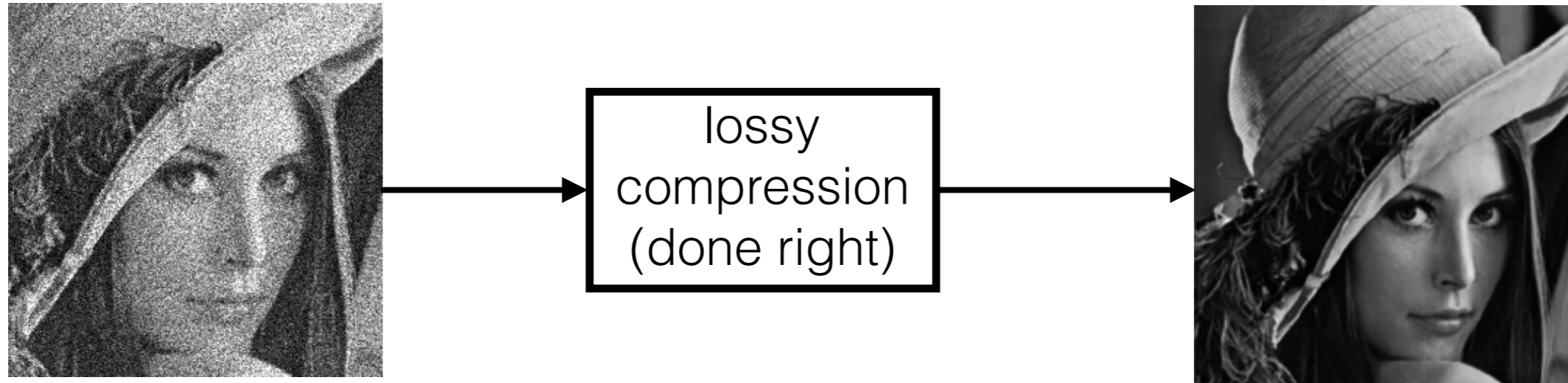- *F-score*: Harmonic mean of sensitivity and precision
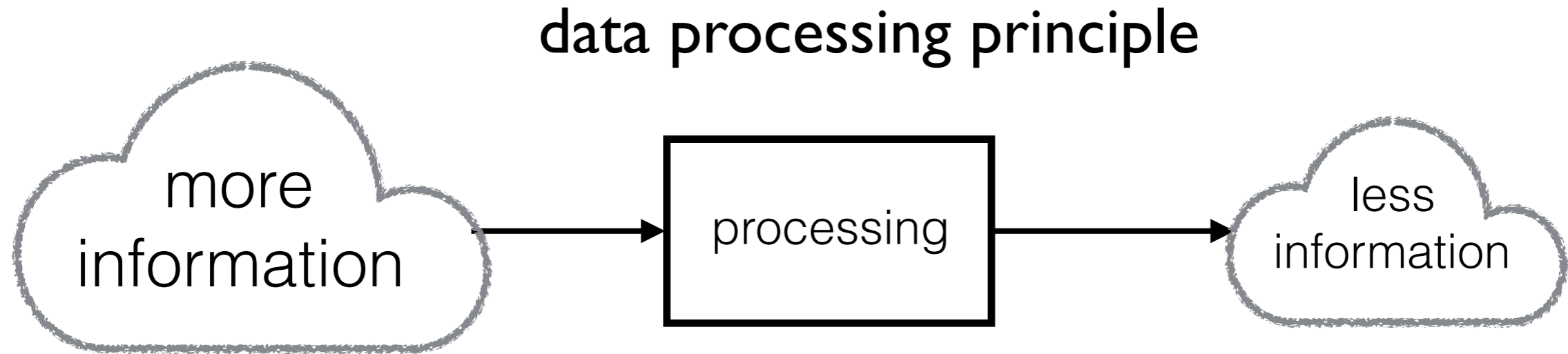
# NA12878, Chr. 11

# zooming in



GATK variant caller

# denoising via lossy compression



- "Occam filters"
- "Compresstimation"
- etc.

**violating the data processing principle ?**

data processing principle



more
information

processing

less
information

suggests need to lift hood of:

- variant callers
- assemblers
-  etc.

# thanks