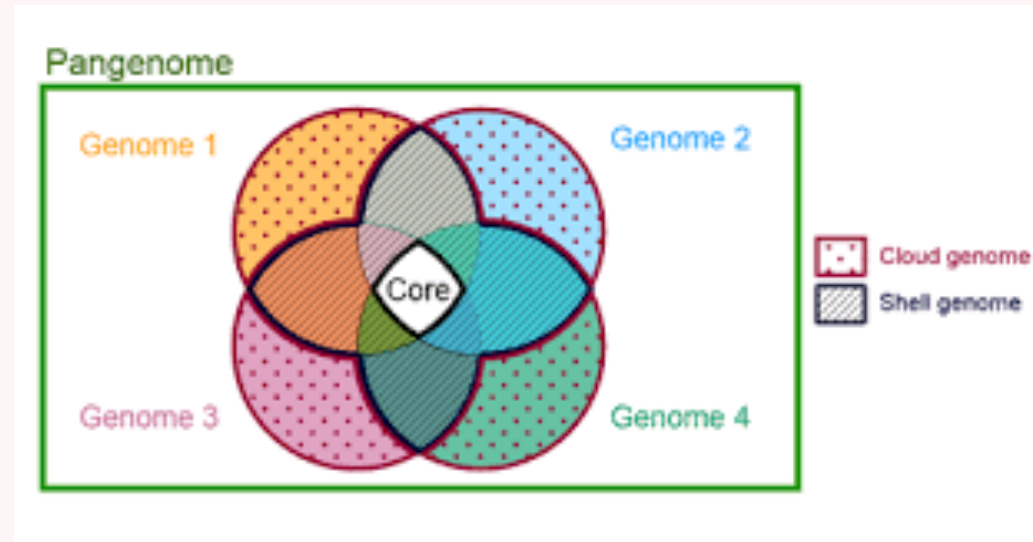# Pangenome Structure?

Ryan C. Fink

# Let's have some fun first...

- Divide in groups of three

- Get a piece of paper and draw a line to get three columns

- Title the three columns "Essential", "Optional", and "Unique"

- Think about your car, your parent's car, or your dream car

- ESSENTIAL: ALL cars must have them to be functioning vehicles

- OPTIONAL: not necessary, but you would love to have them in your car

- UNIQUE: they make a car different from any other (good or bad)
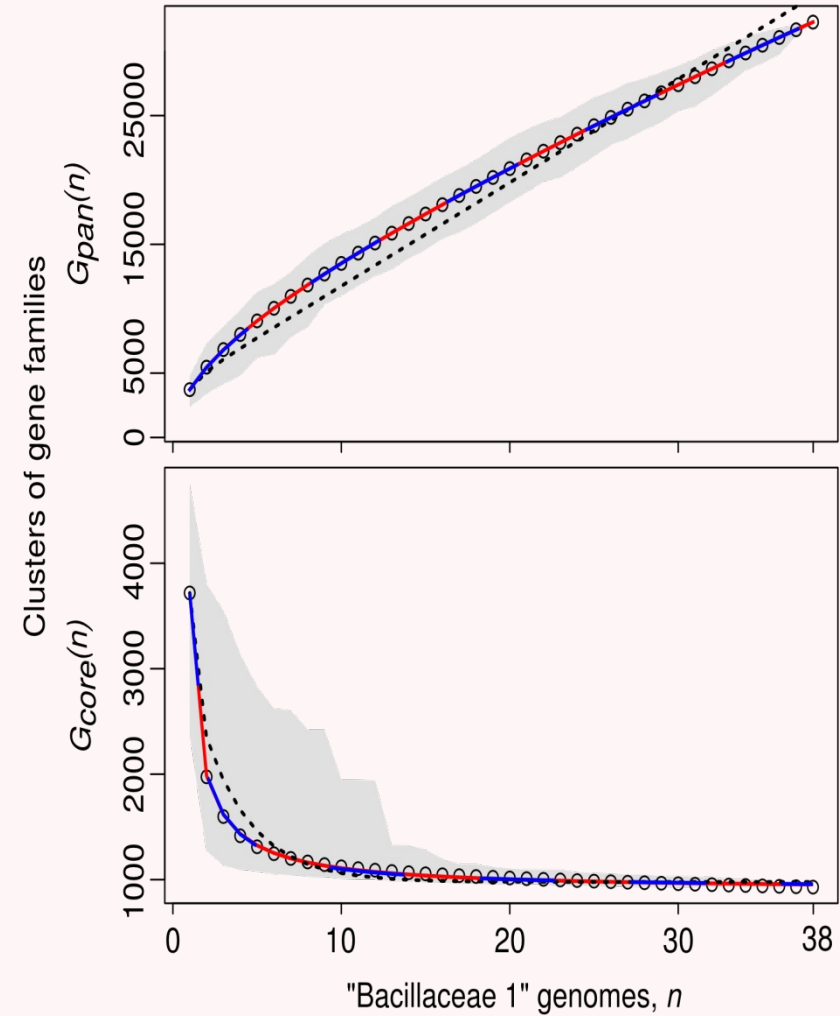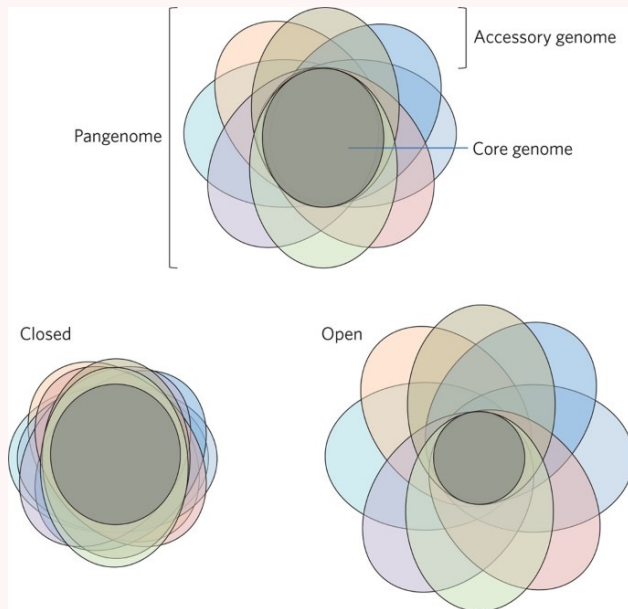
# Pangenome structure

- **PANGENOME** = ESSENTIAL + OPTIONAL + UNIQUE

  - Πας, πασα, παν = all, every, each

- **OPTIONAL** = ACCESSORY (SHELL + CLOUD)

- **UNIQUE** = CLOUD

# Open vs. Close

**CLOSED:** The pangenome size tends to a maximum as number of genomes increases

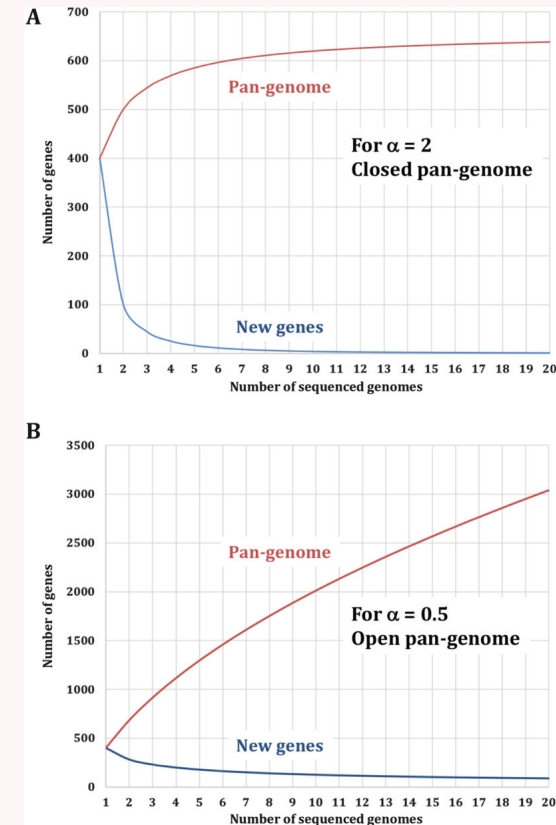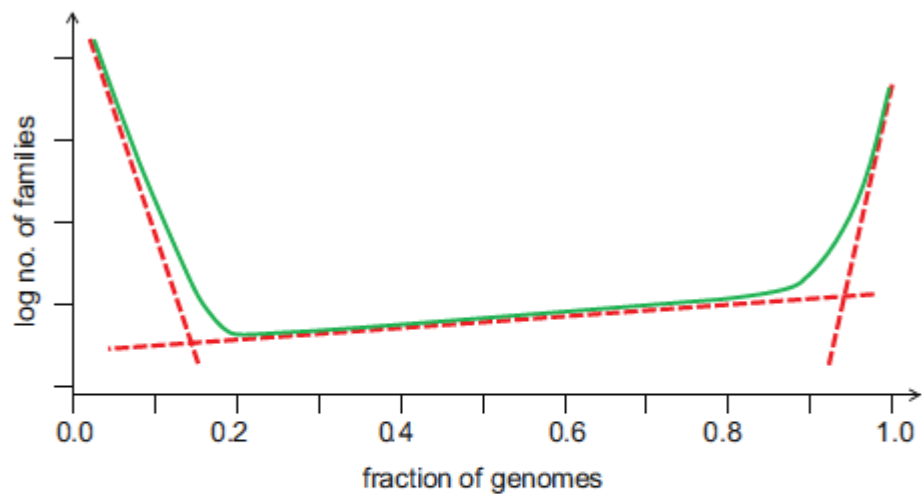**OPEN:** The pangenome keeps increasing as you add new genomes

# Open or closed?

Pangenome measurements follow Heap's law

*"As more and more books are read, the number of different words grows as a power law of the total number of books read"*
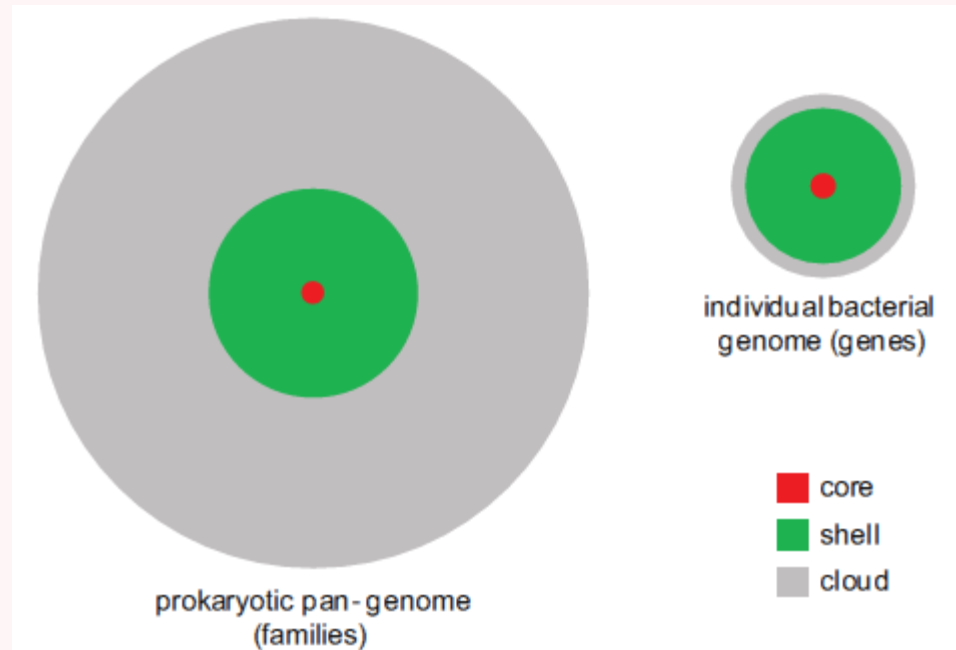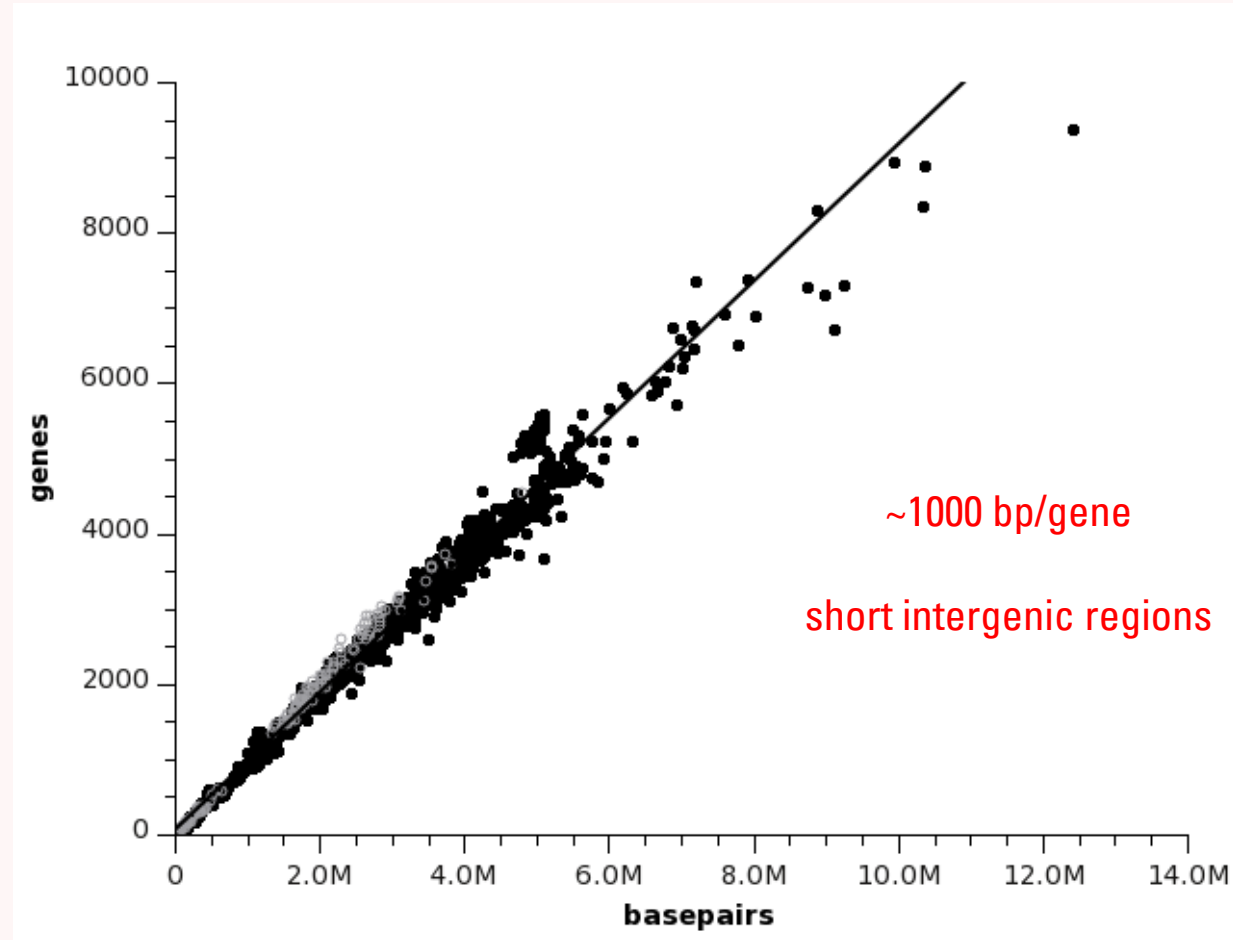
$$n = \kappa N^{-\alpha}$$

# Core, Shell, Cloud… 🤪



FIGURE 4 | The universal distribution of gene commonality in the microbial genomic universe: a generalized schematic. The three broken lines represent three exponential functions that fit the core (on the right), the shell (in the middle) and the cloud (on the left) of prokaryotic genes (O'Malley and Koonin, 2011).
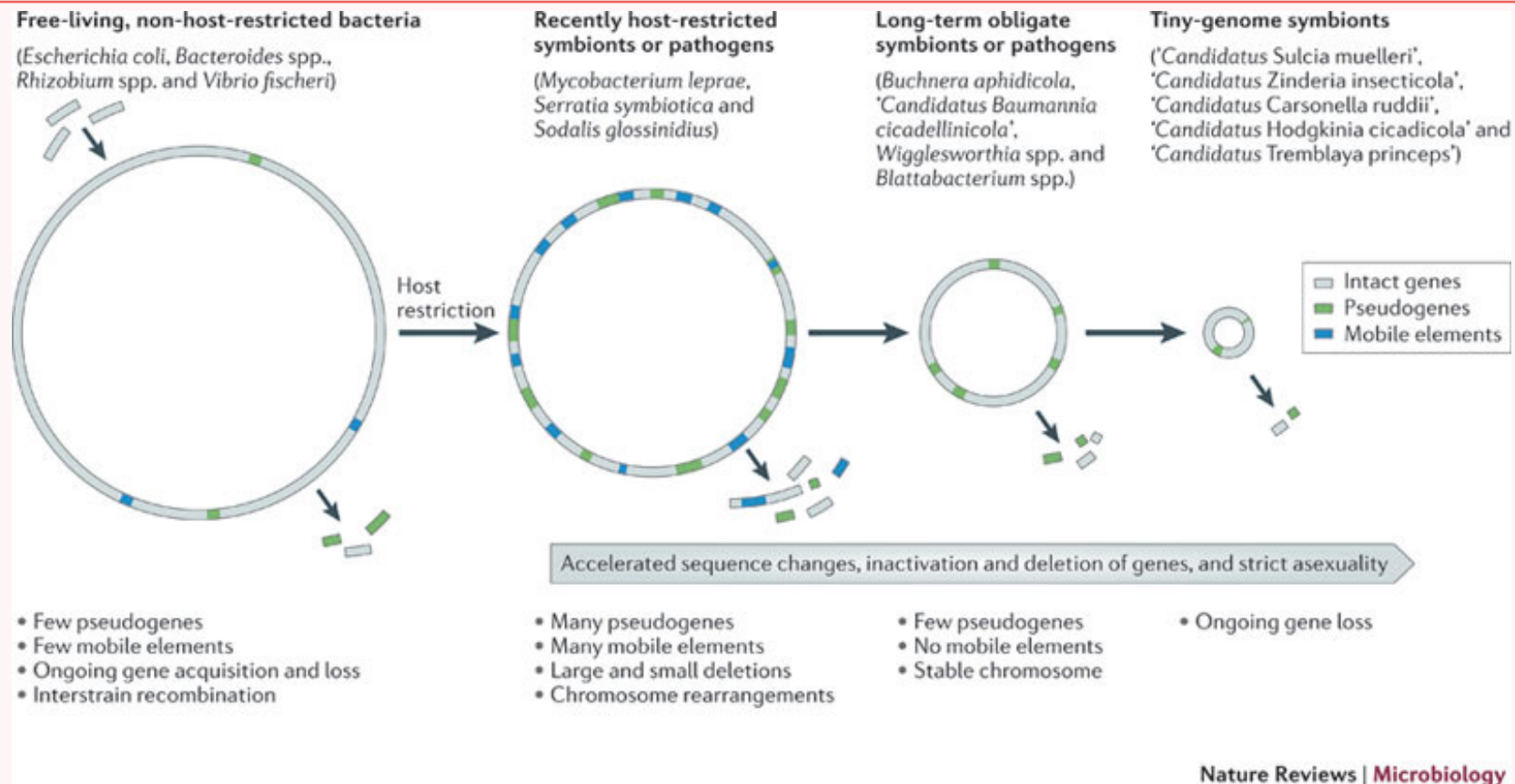


FIGURE 5 | The core, shell, and cloud of microbial genes. A generalized schematic showing the approximate contributions of the core, shell, and cloud to the pangenomes of prokaryotes and individual genomes.

# Genomic streamlining in prokaryotes: a needed little tangent



~1000 bp/gene

short intergenic regions

# It will make sense by the end of this slide…



**Free-living, non-host-restricted bacteria**

(*Escherichia coli, Bacteroides* spp., *Rhizobium* spp. and *Vibrio fischeri*)

**Recently host-restricted symbionts or pathogens**

(*Mycobacterium leprae, Serratia symbiotica* and *Sodalis glossinidius*)

**Long-term obligate symbionts or pathogens**

(*Buchnera aphidicola,* 'Candidatus Baumannia cicadellinicola', *Wigglesworthia* spp. and *Blattabacterium* spp.)

**Tiny-genome symbionts**

('*Candidatus Sulcia muelleri*', '*Candidatus Zinderia insecticola*', '*Candidatus Carsonella ruddii*', '*Candidatus Hodgkinia cicadicola*' and '*Candidatus Tremblaya princeps*')

Host restriction

Intact genes
Pseudogenes
Mobile elements

Accelerated sequence changes, inactivation and deletion of genes, and strict asexuality

- Few pseudogenes
- Few mobile elements
- Ongoing gene acquisition and loss
- Interstrain recombination

- Many pseudogenes
- Many mobile elements
- Large and small deletions
- Chromosome rearrangements

- Few pseudogenes
- No mobile elements
- Stable chromosome

- Ongoing gene loss

Nature Reviews | **Microbiology**

# Horizontal Gene Transfer and Genome Stability

**Horizontal gene transfer:** transfer of genetic information between organisms, as opposed to vertical inheritance from parental organism(s)
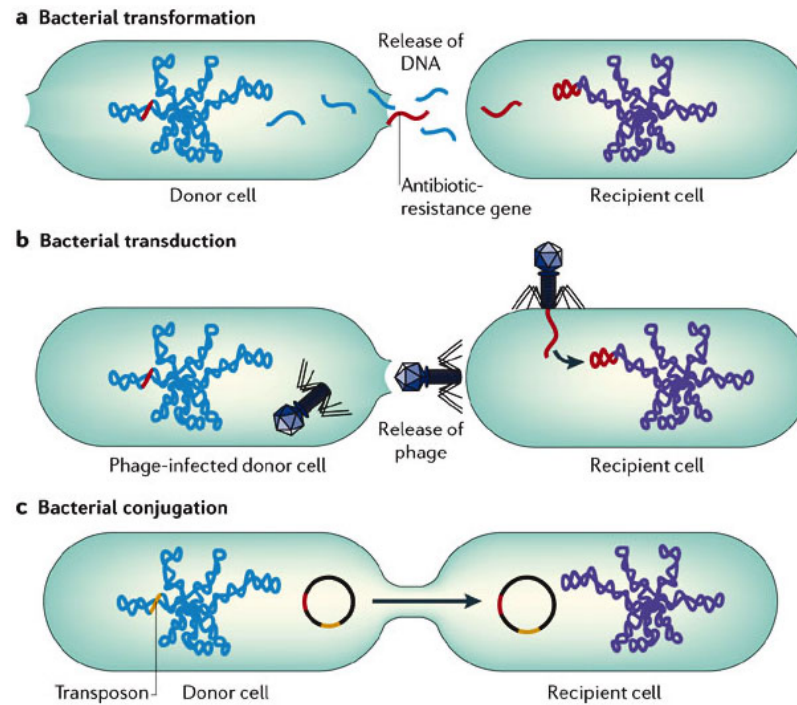
- May be extensive in nature

- May cross phylogenetic domain boundaries

**Detecting horizontal gene flow**

- Presence of genes typically found only in distantly related species

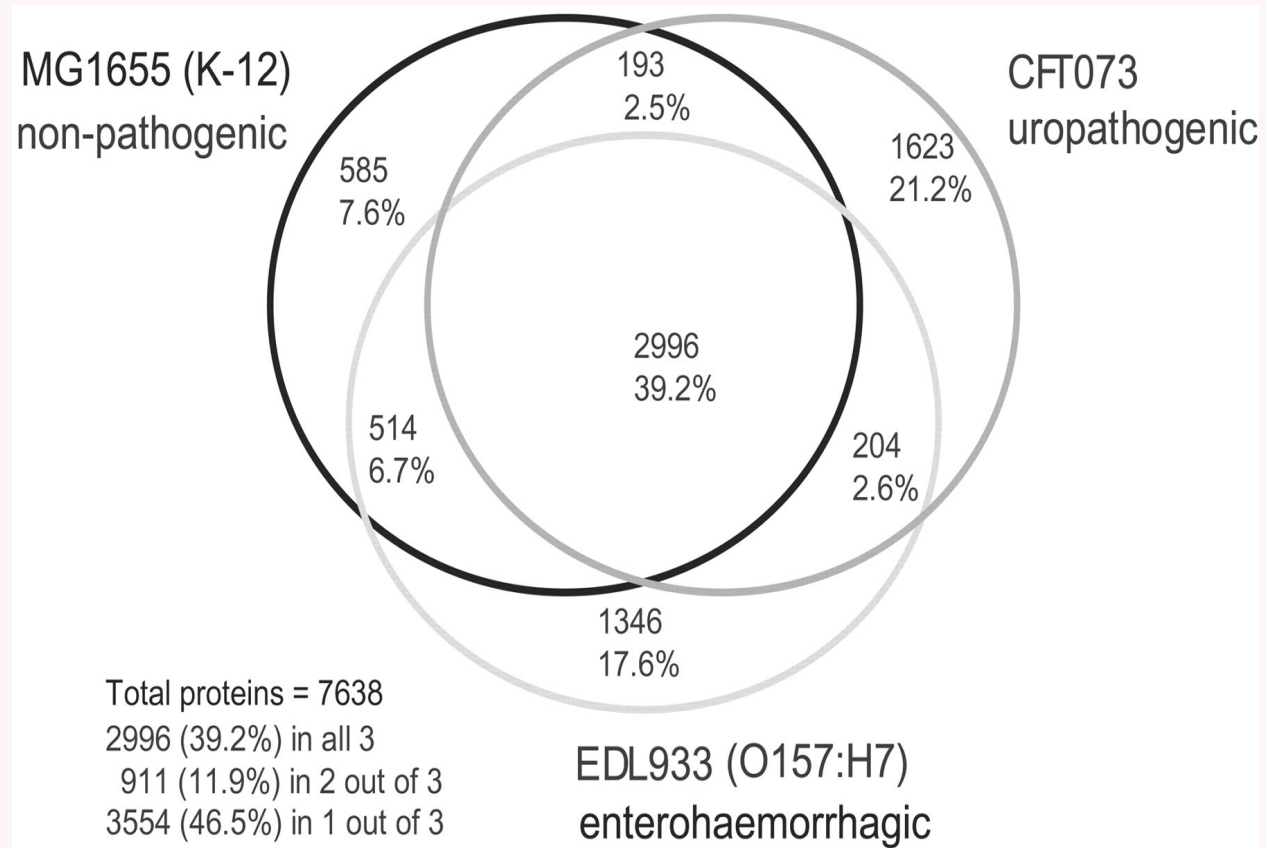- Presence of a DNA with GC content or codon bias that differs significantly from remainder of genome

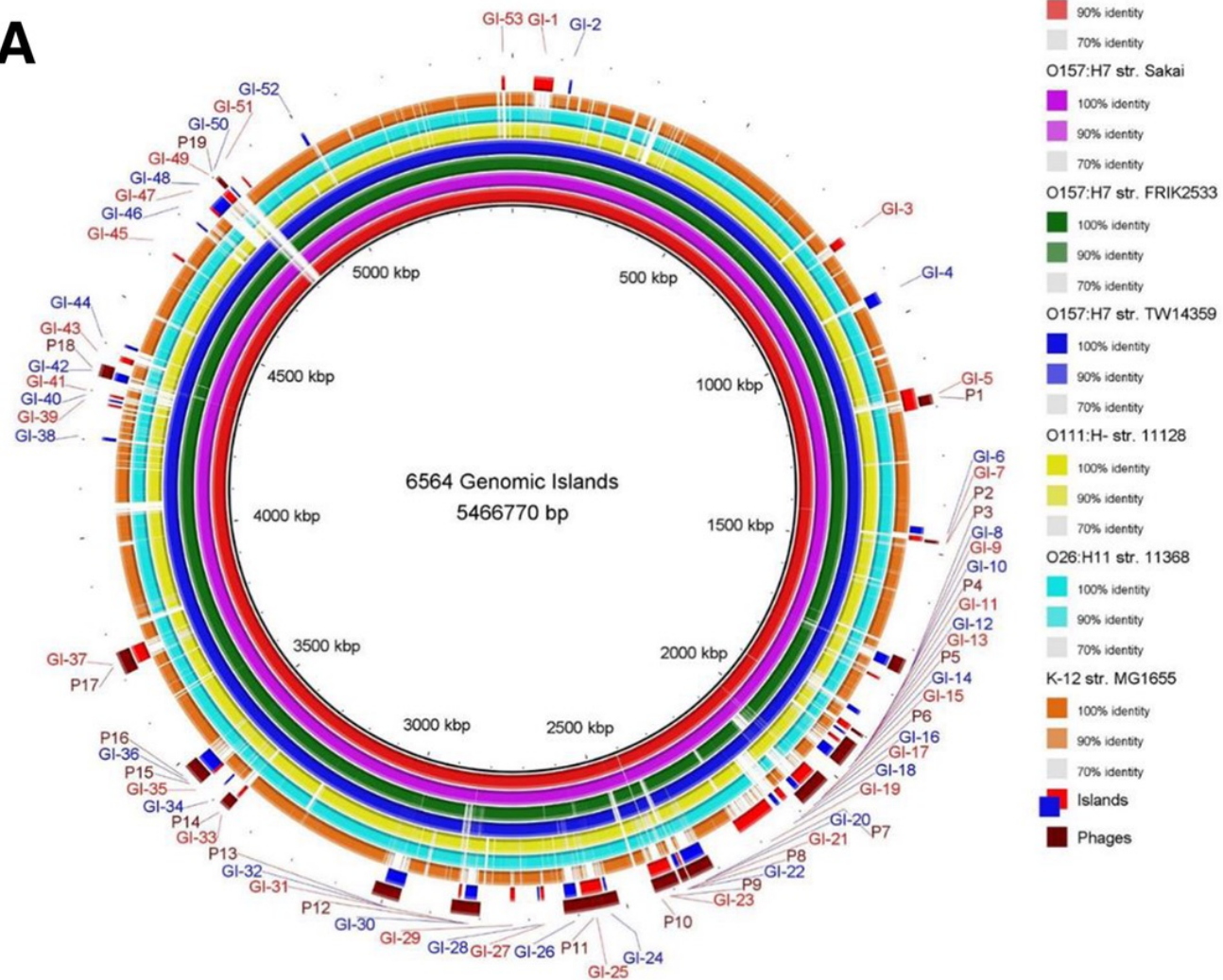**Horizontally transferred genes typically do not encode core metabolic functions**

# Where are the genes coming from?



a **Bacterial transformation**

Release of DNA

Donor cell | Antibiotic-resistance gene | Recipient cell

b **Bacterial transduction**

Release of phage

Phage-infected donor cell | Recipient cell

c **Bacterial conjugation**

Transposon | Donor cell | Recipient cell

Copyright © 2006 Nature Publishing Group
**Nature Reviews | Microbiology**

# Gene Content Variation among *E. coli* genomes. Evidence for horizontal transfer –



MG1655 (K-12) non-pathogenic

CFT073 uropathogenic

193 2.5%

585 7.6%

1623 21.2%

2996 39.2%

514 6.7%

204 2.6%

1346 17.6%

Total proteins = 7638
2996 (39.2%) in all 3
911 (11.9%) in 2 out of 3
3554 (46.5%) in 1 out of 3

EDL933 (O157:H7) enterohaemorrhagic

Welch et al (2002).

# Accessory Genome: a treasure trove of information

- Host specificity

- Lifestyle

- Potential vectors

- Superpowers!

# How can we find the accessory, juicy genes?

## Different softwares

- **Roary**
  - 98% ID, 98% match length
  - No easy way to change them

- **Panaroo**
  - 98% ID, 98% match length
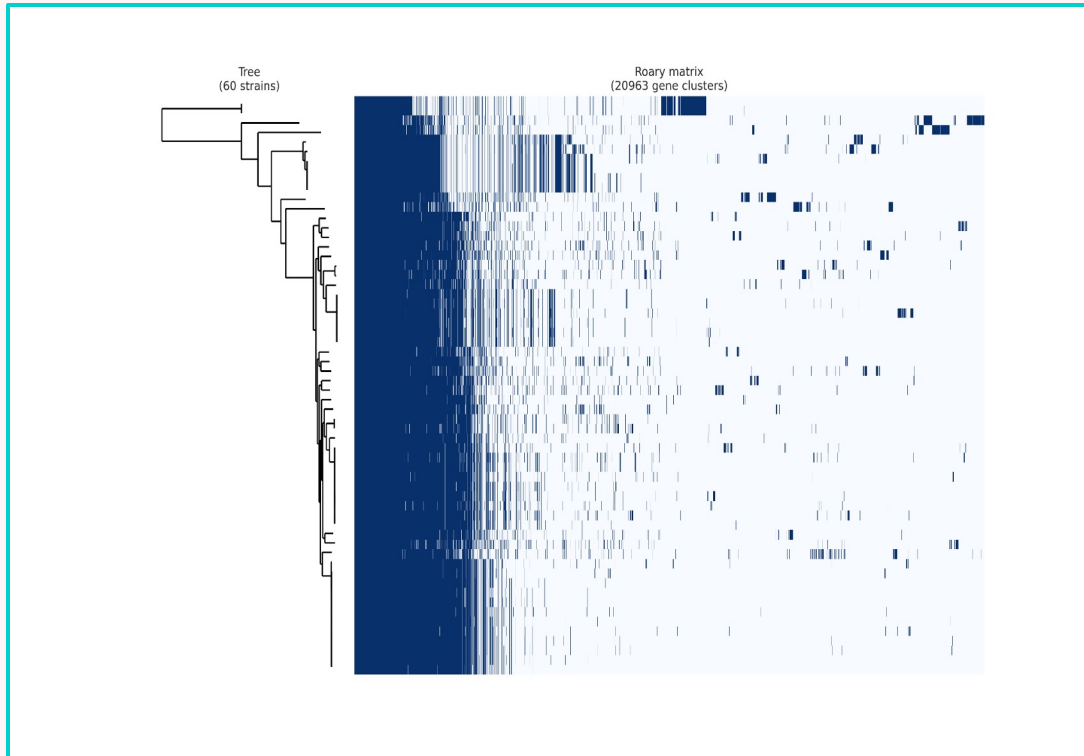  - Easy way to change them

- **PPanGGOLiN**
  - It's… different



**Figure S1:** A comparison of the pipelines used by different pangenome analysis tools.

# Pet Set of *Salmonella*

# ROARY

- ASSUMPTION:
  - All entries from the same species
- Filter and precluster the proteins
- All against all comparison using BLASTP
- Cluster with MCL



Pan-Genome Construction

16

# Roary's Output

# Panaroo

- Panaroo corrects annotation errors

- Contamination appears in the graph as poorly supported components

- Genes are often mis-annotated near contig breaks

- Corrects same DNA sequence translated in multiple reading frames

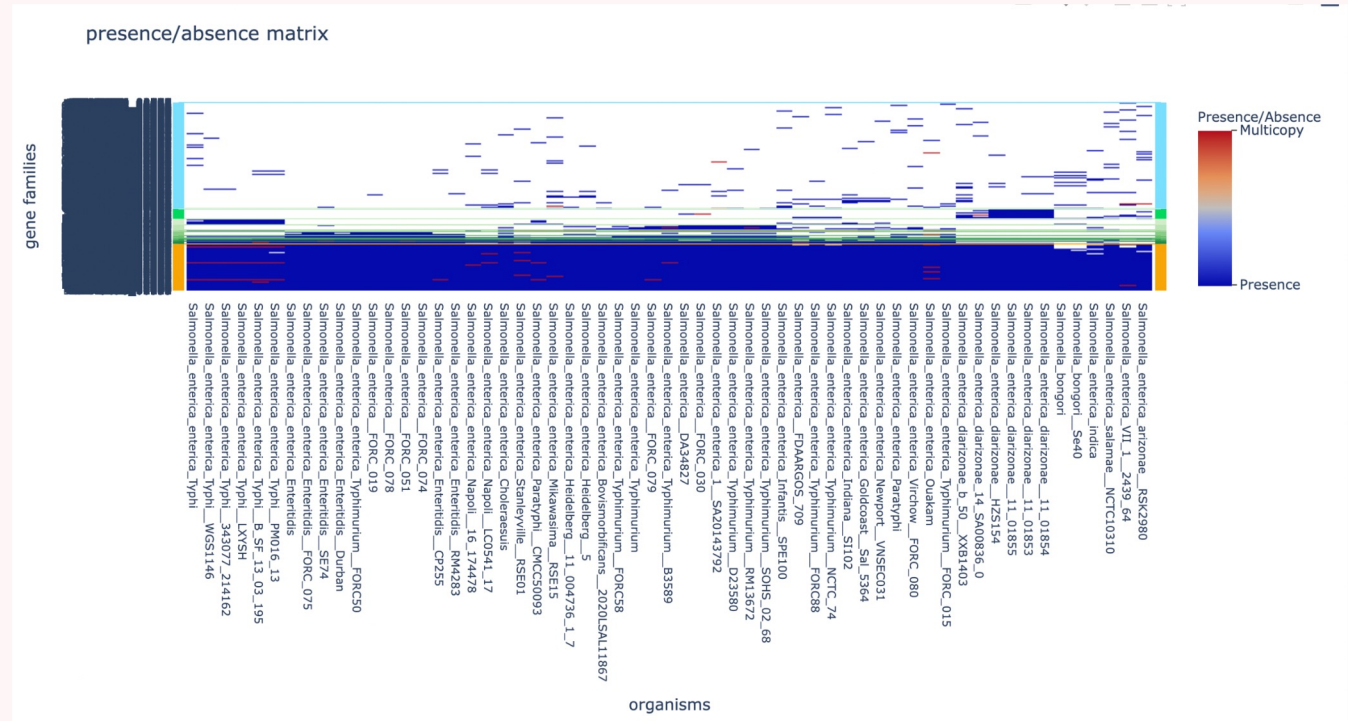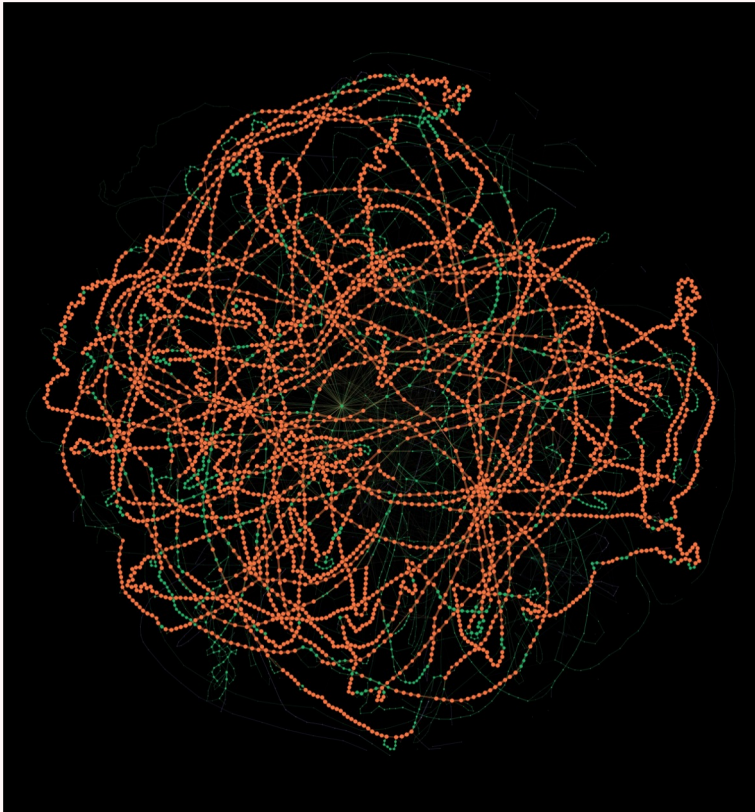- Combine diverse gene families into a single gene

# Panaroo's Output

# PPanGGOLiN

- Clustering through MMSeqs

- Statistical analysis to assign families to partitions

- No reference tree (but can be produced)

- Can be used for different taxonomic levels

- PPanGGOLiN returns a partitioned pangenome graph

    - persistent, shell and cloud partitions
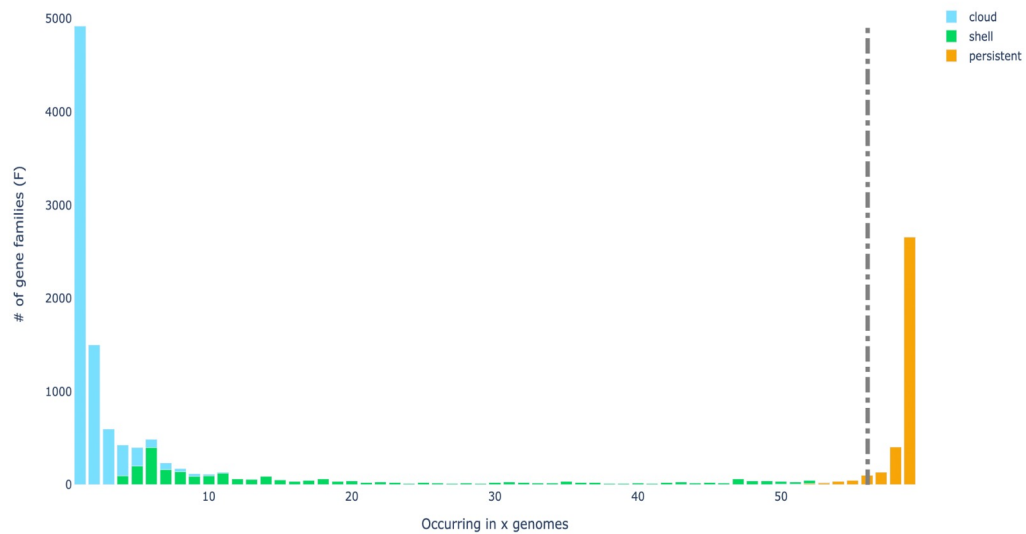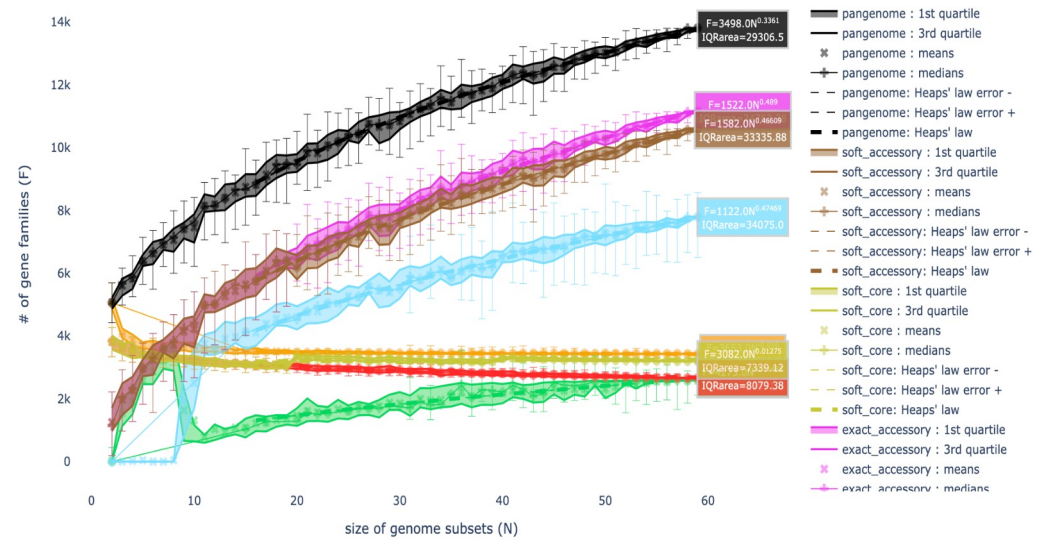    - overlaid on the neighborhood graph

# PPanGGOLiN's Output

# PPanGGOLiN's Output -2

# How to Choose a Tool?

**Considerations:**

- Desired visualization(s)

- Data set size and computational resources

- Sensitivity vs specificity