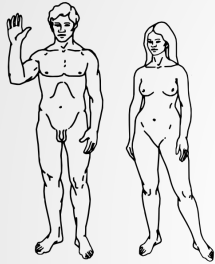


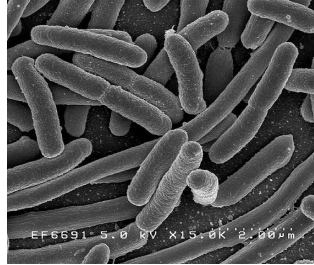
Infectome & Novel Bug Pipeline

Mark Kowarsky, Quake Lab

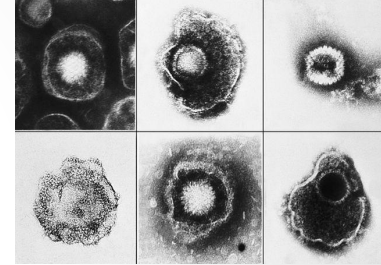
“Human” DNA



- 10^{13} cells
- Gb genome
- 10^{22} nucleotides



- 10^{14} cells
- Mb genome
- 10^{20} nucleotides



- 10^{14} virions
- kb genome
- 10^{17} nucleotides

- ~1% non-human, sequence cell-free DNA from blood

Pipeline overview

Samples

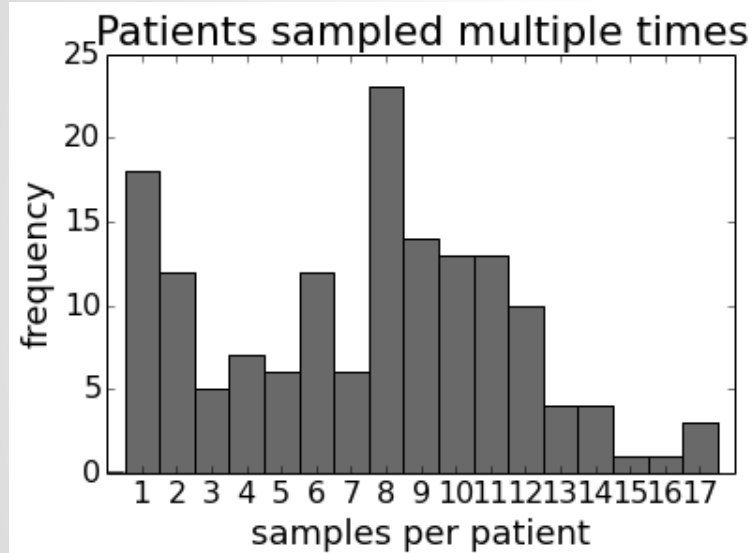
1. Preprocessing
2. Host removal
3. Infectome [known]
4. Finding the gold
[unknown]

Sets of samples

1. Aggregated assembly
2. Realign for coverage
3. Identify contigs

Samples

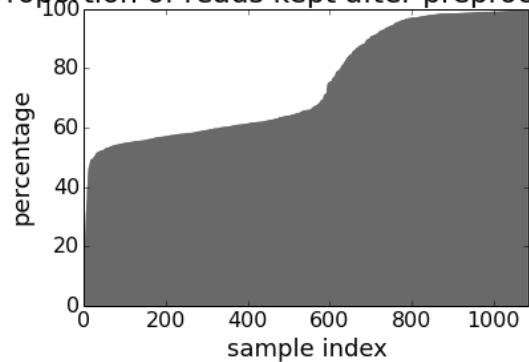
- 400+ people
- 1000+ samples
 - 642 heart transplant
 - 389 lung transplant
 - 73 bone marrow transplant
 - 18 chronic fatigue syndrome
- Illumina sequenced, 2x100bp or 1x50bp



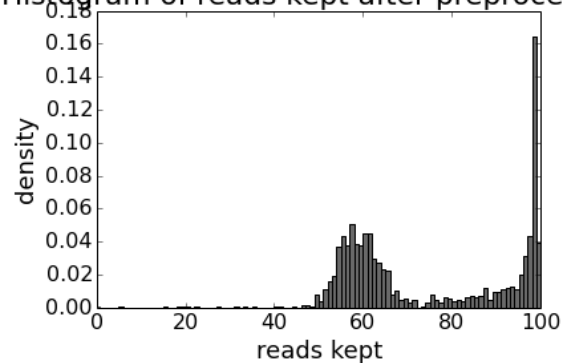
1. Preprocess

- Quality check (fastQC)
- Trim adapters
- Merge overlapping reads
- Set of 'good' reads to use

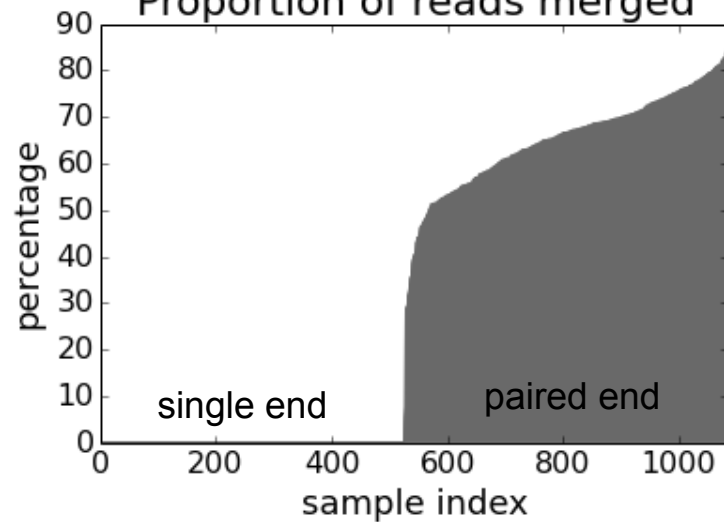
Proportion of reads kept after preprocessing



Histogram of reads kept after preprocessing



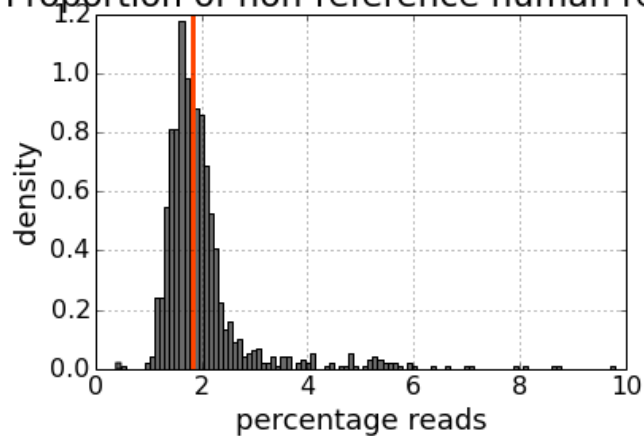
Proportion of reads merged



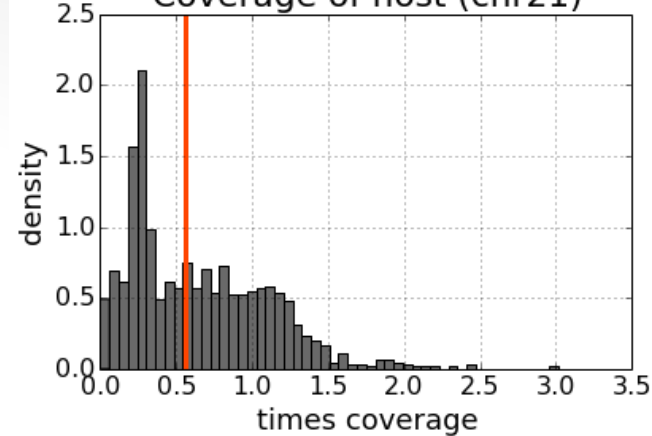
2. Host removal

- Align to hg19
- Align to phiX
- Remove low complexity reads
- ~2% of reads remain

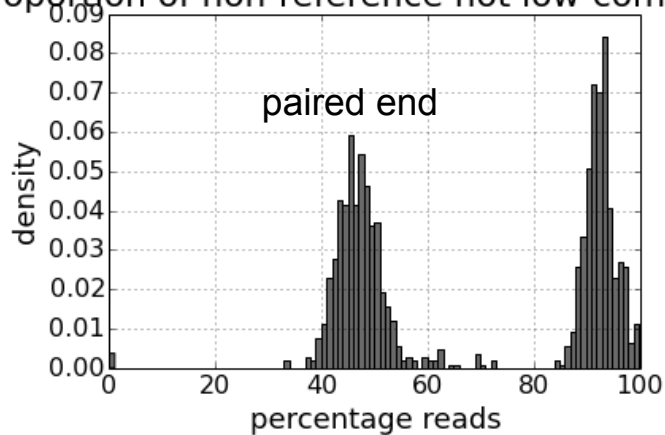
Proportion of non-reference human reads



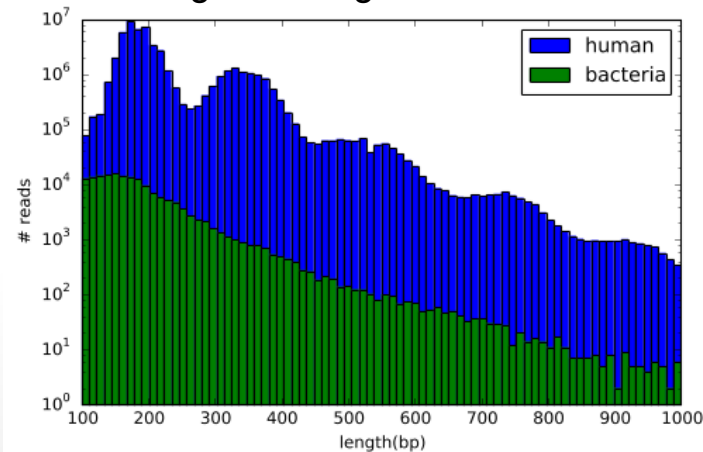
Coverage of host (chr21)



Proportion of non-reference not low-complexity



Fragment length differences



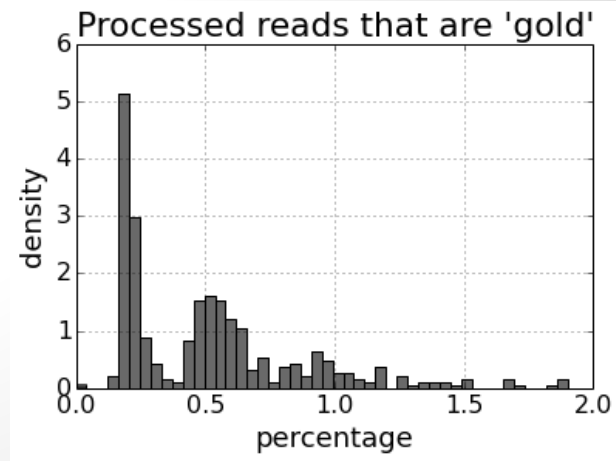
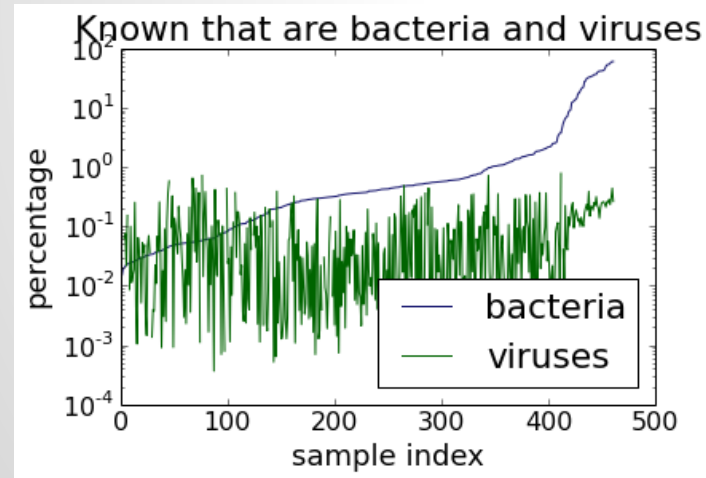
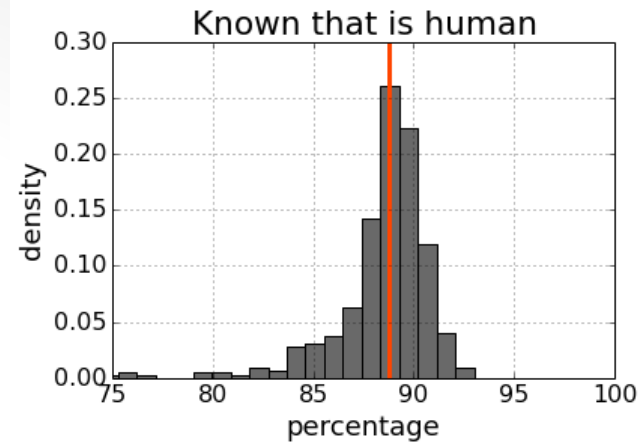
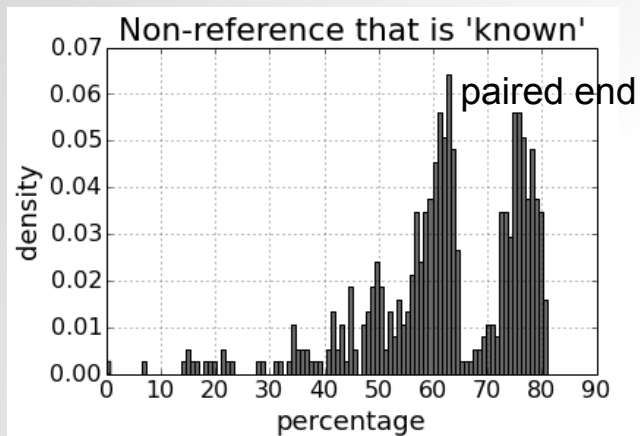
3. Known infectome

- BLAST filtered non-reference reads against bug database
- Use GRAMMy* to estimate abundance
- See Lance's talk for details

* Genome Relative Abundance using Mixture Models - <https://bitbucket.org/charade/grammy/wiki/Home>

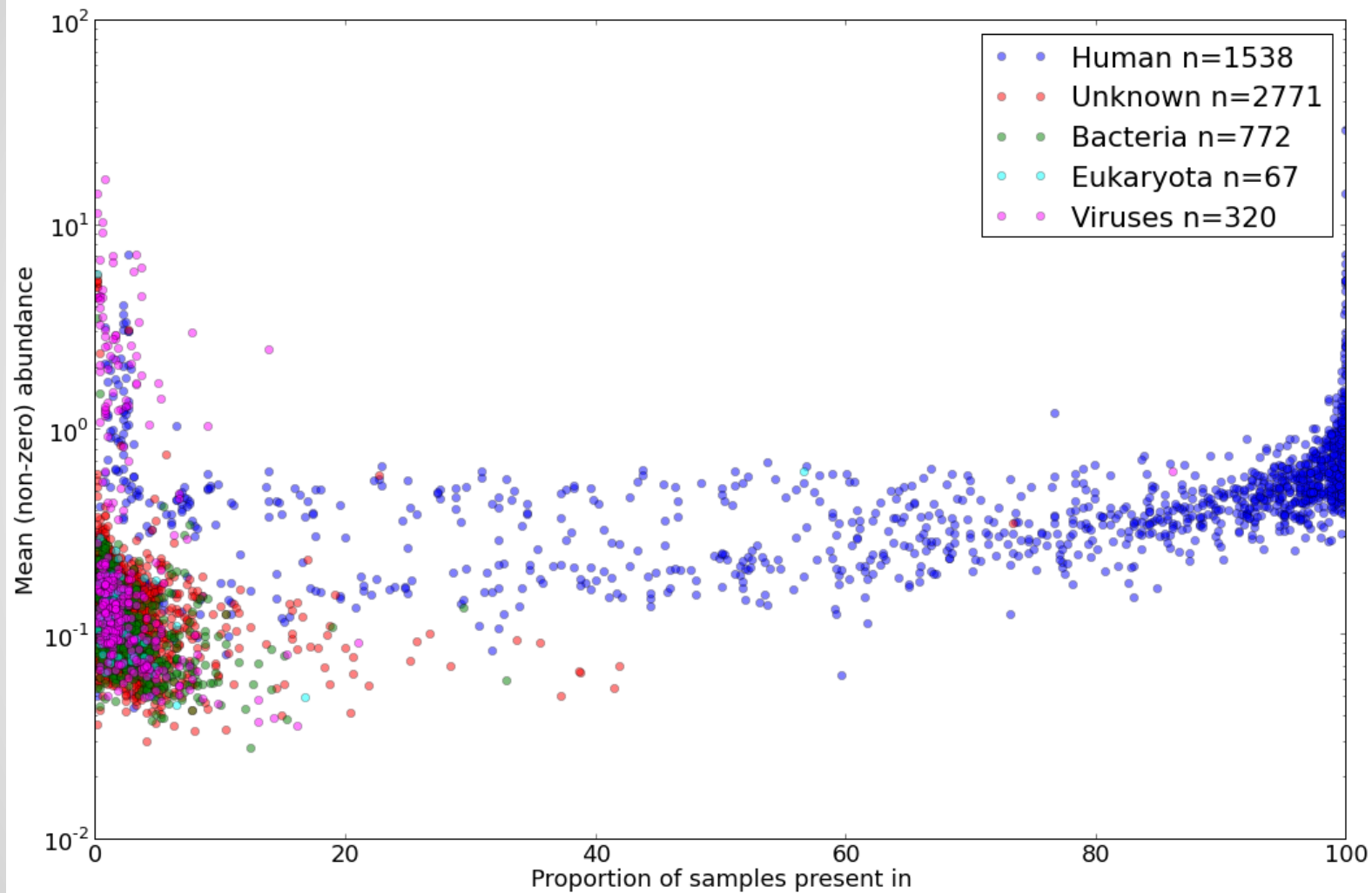
4. Novel bugs - “gold”

- BLAST against NCBI's NT database
- Filter reads that have excellent alignments
- Remnants are 'gold'



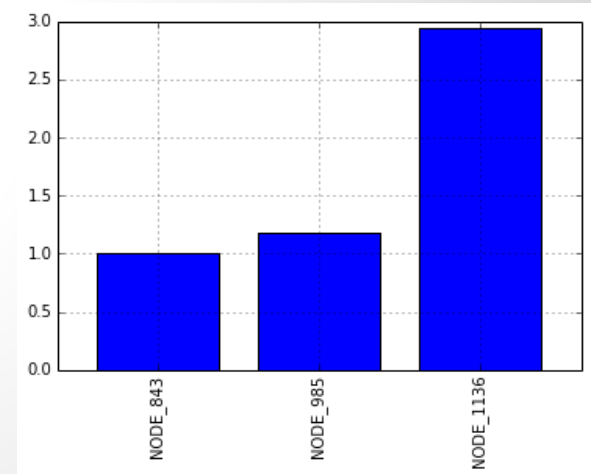
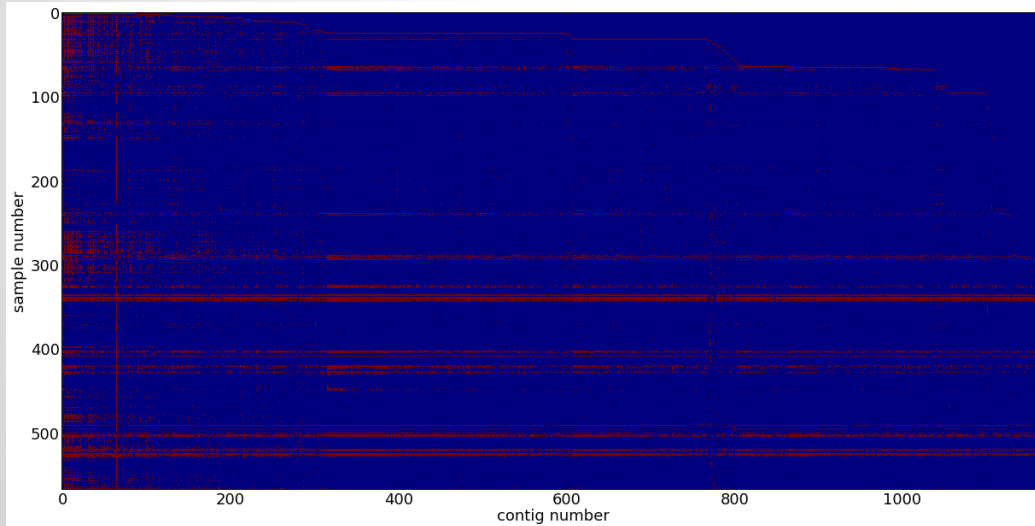
Sets

- Aggregate reads from groups of samples, e.g.
 - same transplant type
 - same patient
 - hospital
- Assemble
- Realign
- Identify



Clustering

- Correlate novel contigs
- Find clusters highly associated



Further work

- Finish running pipeline on data (~1-2 weeks)
- Find good bug candidates
 - high coverage
 - clinical correlation
- PCR to validate assembly

Acknowledgements

- Mickey Kertesz
- Lance Martin
- Iwjin De Vlaminck