

PATRIC Bioinformatics Resource Center

Genome Assembly in PATRIC

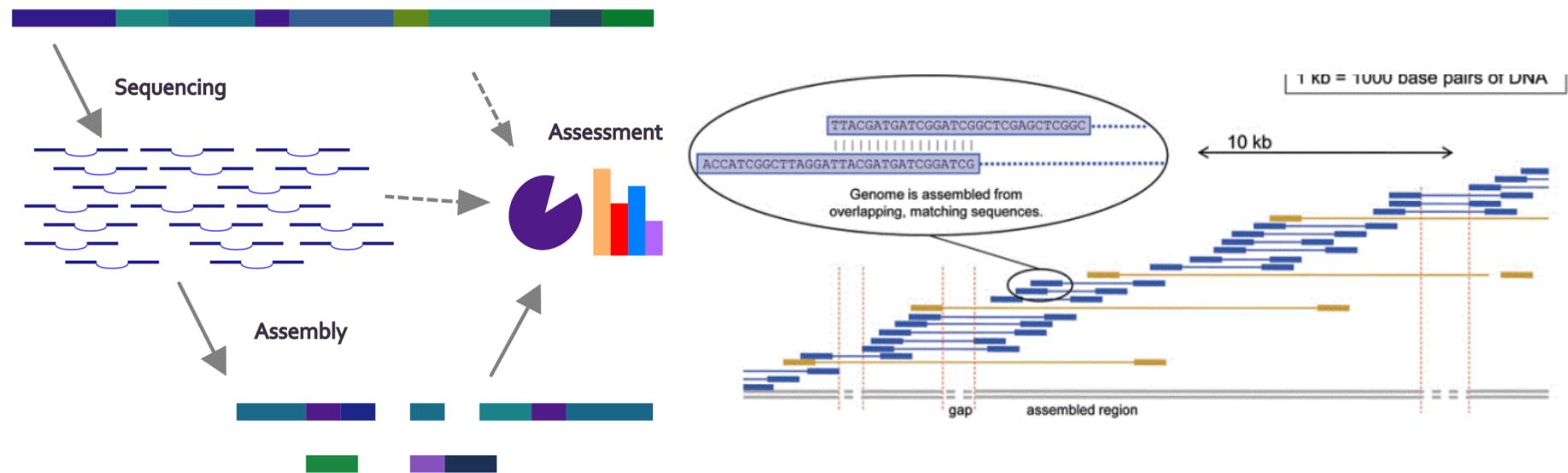
Presented by Neal Conrad

Slides originally by Fangfang Xia



The Sequence Assembly problem

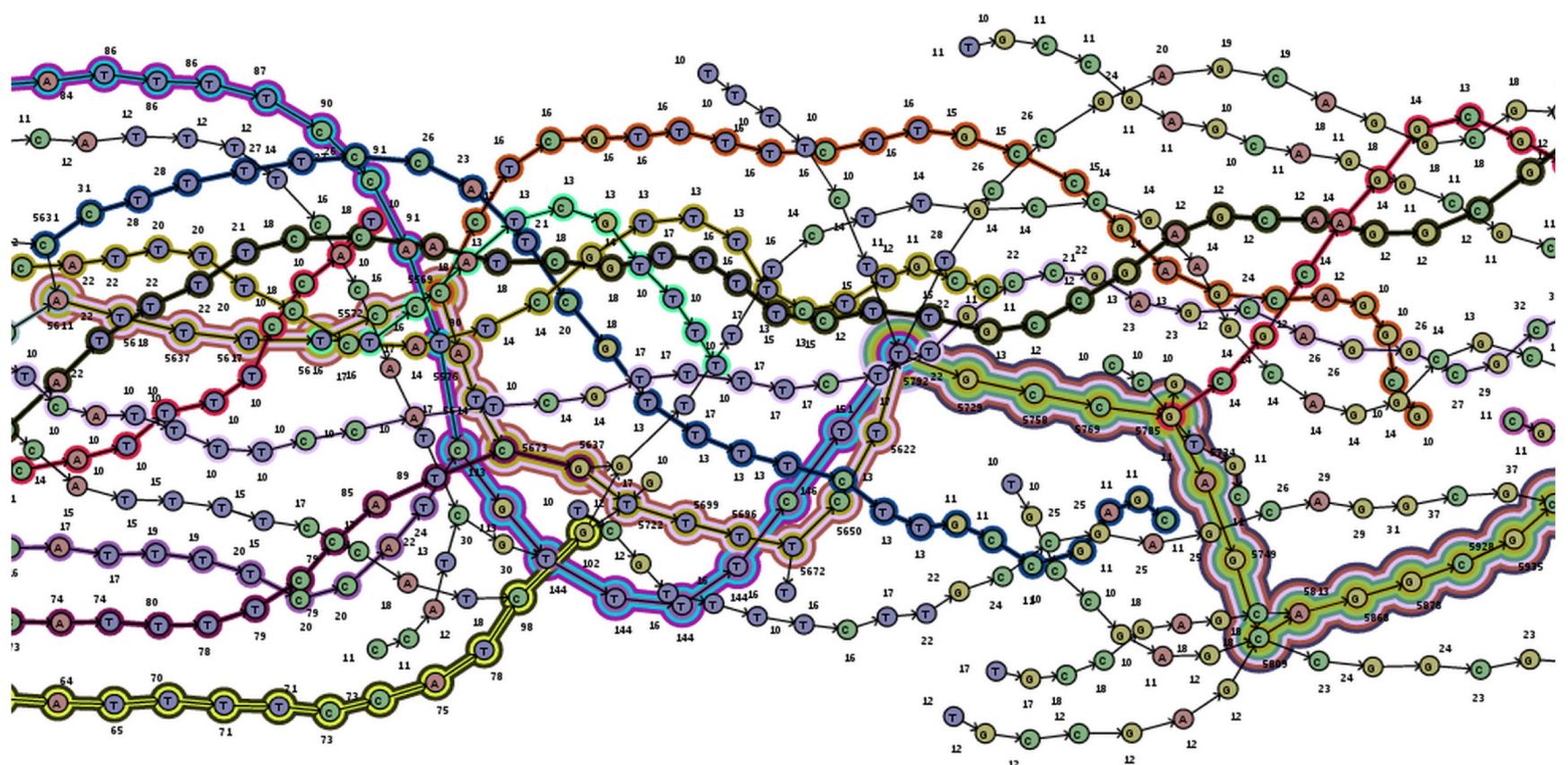
- ▶ Reconstructing contiguous DNA regions (contigs) from a set of short sequences (reads)



An incomplete list of assemblers

Abyss, AllPaths, AllPaths-LG, AMOS, Arapan-M,
Arapan-S, Celera, CLC, Clustgun, Cortex,
Discovar, DNA Baser, Dragon, Edena, Euler,
Euler-sr, FERMI, Forge, Geneious, Graph
Constructor, HGAP, IDBA, IDBA-UD, Kiki, Meta-
velvet, Minia, MIRA, NextGENe, Newbler,
PADENA, PASHA, Phrap, Ray, Ray-meta, REAPR,
Sequencher, SeqMan, SGA, SHARCGS, SOPRA,
SSAKE, SOAPdenovo, SPAdes, Staden, Taipan,
TIGR, VCAKE, Phusion, QSRA, Velvet, YAGA

Assembly graph

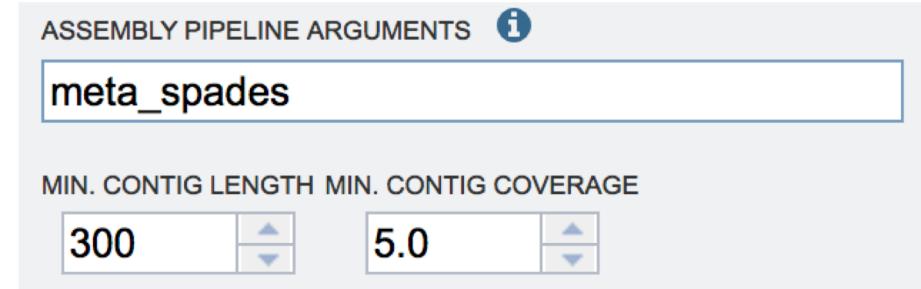
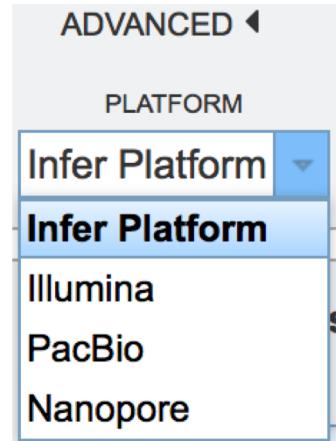


"Assemblers should take in your data and automatically do the best possible job with it."

– A reviewer for Assemblathon

Common Assembly Scenarios

- ▶ *Short read assembly*— Illumina sequencers
- ▶ *Long read assembly*— PacBio, Oxford Nanopore
- ▶ *Hybrid assembly*— Short reads + long reads
 - Short for assembly + long for scaffolding (spades)
 - Long for assembly + short for finetuning (auto; planned)
- ▶ *Whole genome assembly*
- ▶ *Plasmid assembly*
- ▶ *Metagenome assembly*
- ▶ *De novo assembly vs reference-guided assembly*
 - lib.consensus.fa output in variation analysis



ASSEMBLY SERVICE

COMPUTE CAPABILITIES

PREPROCESSING	ASSEMBLY	POST-ASSEMBLY	EVALUATION
SolexaQA	Kiki	SSPACE:Scaffold	ALE
Length Filtering	SPAdes	REAPR: Break	REAPR
SGA: Q-trim / Q-filter	SGA	GAM-NGS Merge	QUAST
TagDust: Adapter Removal	MaSuRCA	...	
BayesHammer EC	Discovar		
SGA: Error Correction	IDBA-UD		
...	Velvet	BWA	Nx Plots
	A5	Bowtie2	Contig Length
	ALE Comparison
			Benchmarks

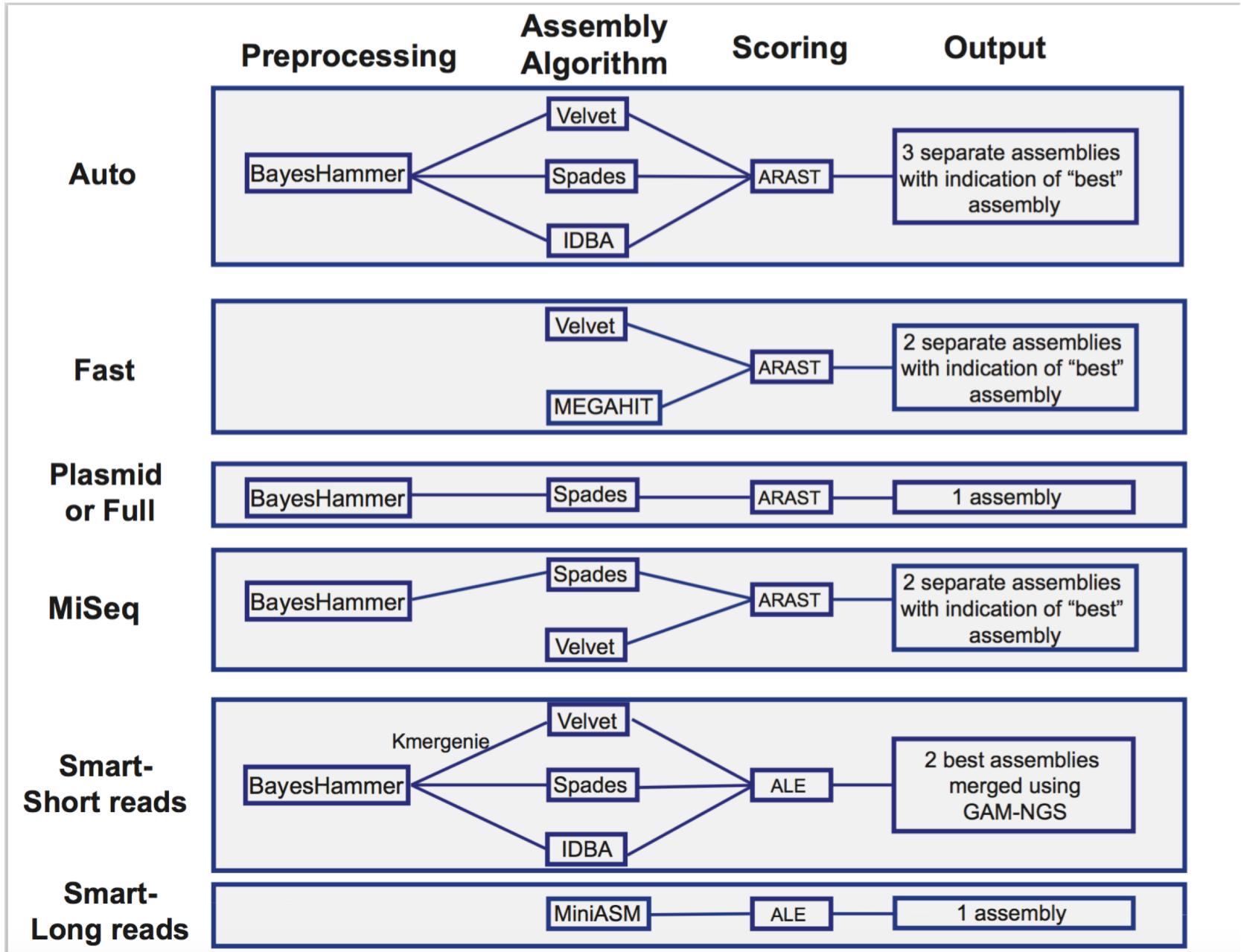
Curated assembly strategies

Parameters 

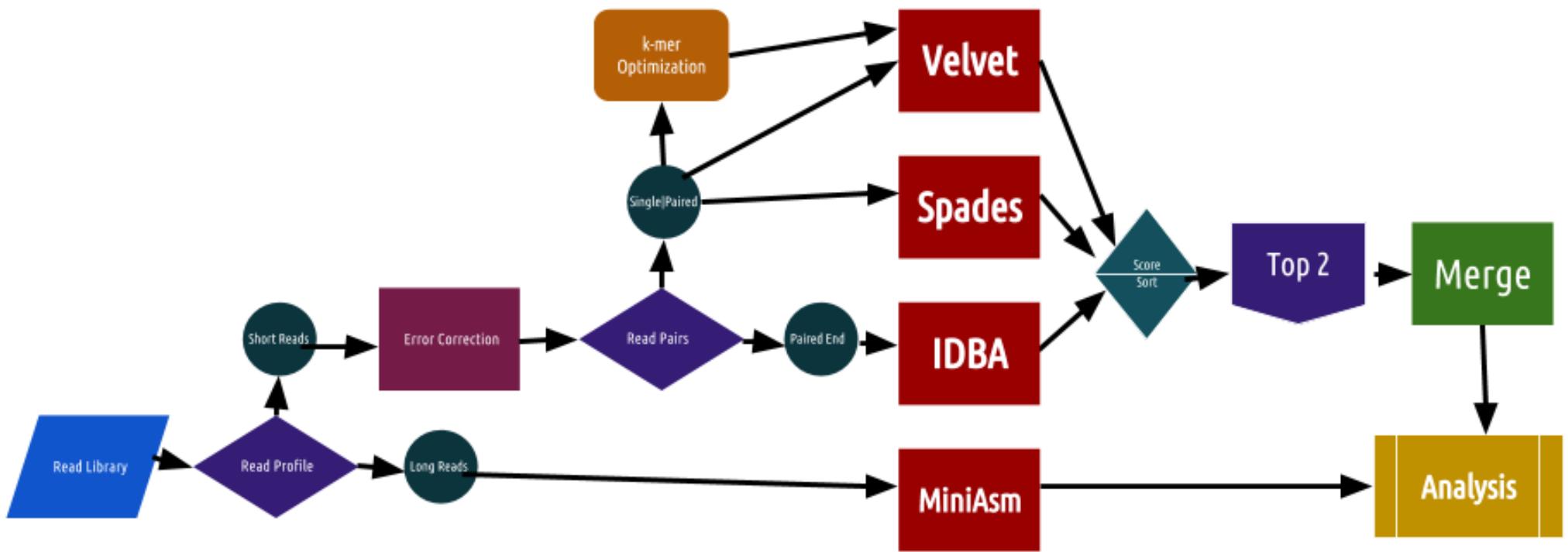
ASSEMBLY STRATEGY

- auto
- auto**
- fast
- full spades
- kiki
- miseq
- plasmid
- smart

ADVANCED ▼



The “smart” assembly recipe



Which assembly recipe to use

- ▶ *auto* — the evolving default strategy recommended for most data
- ▶ *full spades* – runs the full SPAdes pipeline, one of the best assemblers for microbial genomes
- ▶ *fast* — ~2X faster than *auto*; suited for large genomes or simple microbial communities (velvet + megahit)
- ▶ *kiki* — very fast but does not use paired end information; good for metagenome assembly
- ▶ *miseq* — good for Illumina MiSeq reads that are 250–350 bp long (Spades with more k-mer iterations)
- ▶ *smart* — the slowest and sometimes the most accurate
- ▶ *plasmid* — plasmid assembly (plasmidSPAdes)

Typical execution times

for a typical microbial genome

Recipe	Hours
smart	3 ~ 100
auto / miseq	2 ~ 80
fast	1 ~ 12
kiki	1 ~ 6

Depends on read depth, sequencing errors, genome size, repeat structure, etc.

Evaluate assembled contigs

QUAST report

19 June 2014, Thursday, 08:19:51

All statistics are based on contigs of size ≥ 500 bp, unless otherwise noted (e.g., "# contigs (≥ 0 bp)" and "Total length (≥ 0 bp)" include all contigs.)

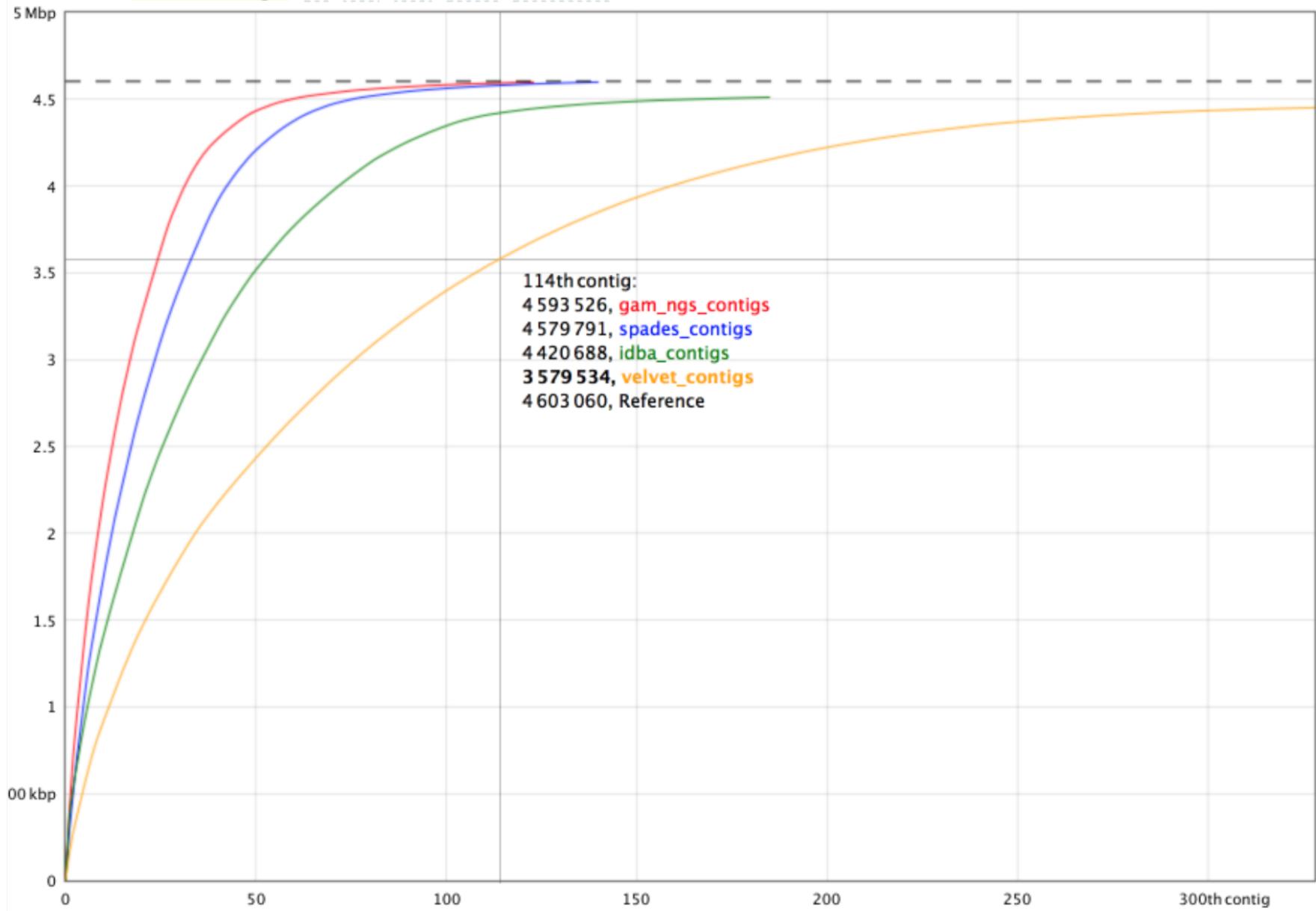
Extended report

worst.....best

Genome: 4 603 060 bp, G+C content: 68.79%

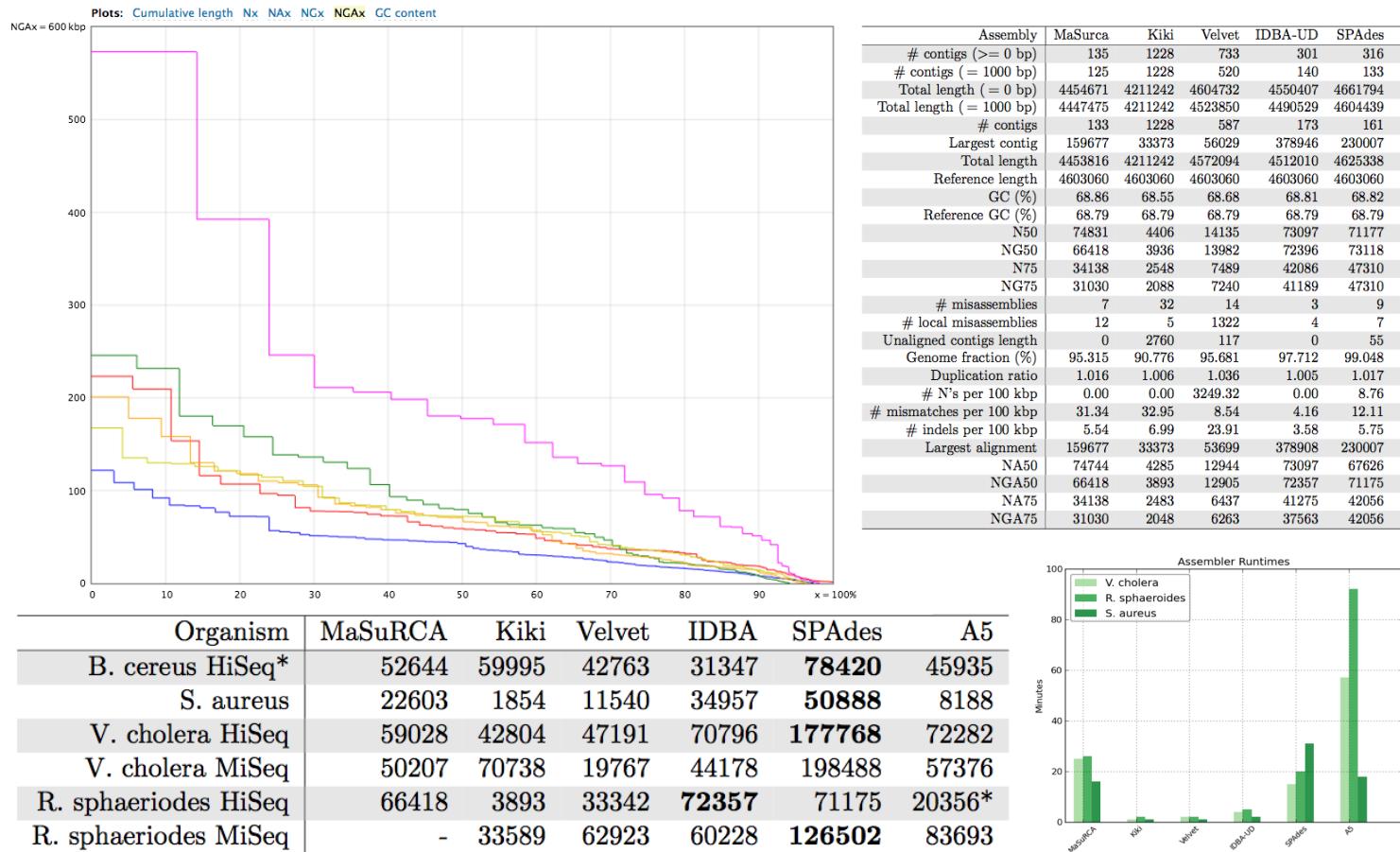
Statistics without reference	gam_ngs_contigs	spades_contigs	idba_contigs	velvet_contigs
# contigs	123	140	185	328
Largest contig	399 330	292 708	378 933	160 339
Total length	4 598 516	4 598 552	4 511 207	4 451 511
N50	130 248	95 827	61 227	26 869
Misassemblies				
# misassemblies	6	6	1	29
Misassembled contigs length	153 182	150 167	24 718	876 583
Mismatches				
# mismatches per 100 kbp	16.080	16.52	6.08	12.69
# indels per 100 kbp	4.020	3.91	3.5	10.210
# N's per 100 kbp	0	0	0	373.58
Genome statistics				
Genome fraction (%)	98.898	98.895	97.953	96.39
Duplication ratio	1.01	1.01	1.001	1.003
NGA50	130 247	95 827	61 162	24 307
Predicted genes				
# predicted genes (unique)	4480	4501	4464	4695
# predicted genes (≥ 0 bp)	4519	4540	4464	4695
# predicted genes (≥ 300 bp)	3970	3984	3919	3989
# predicted genes (≥ 1500 bp)	550	545	537	480
# predicted genes (≥ 3000 bp)	44	44	38	34

Plots: Cumulative length Nx NAX NGx NGAx GC content



AssemblyRAST: Assembler Comparison

```
for LIB in $(ls)
do
    ar_run -f $LIB/rd*.fq -a masurca kiki velvet spades idba a5
done
```



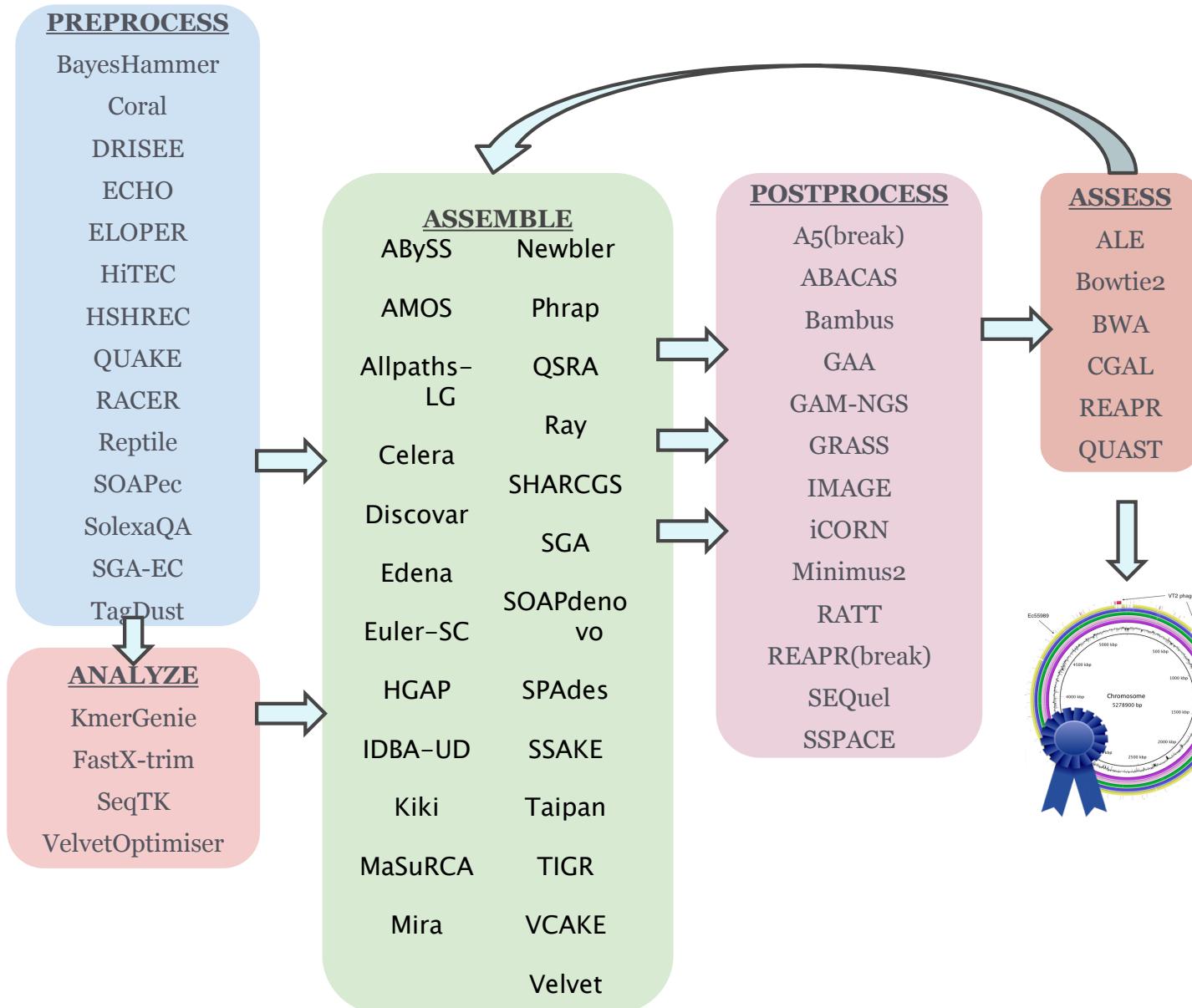
Interactive Demo

Acknowledgements

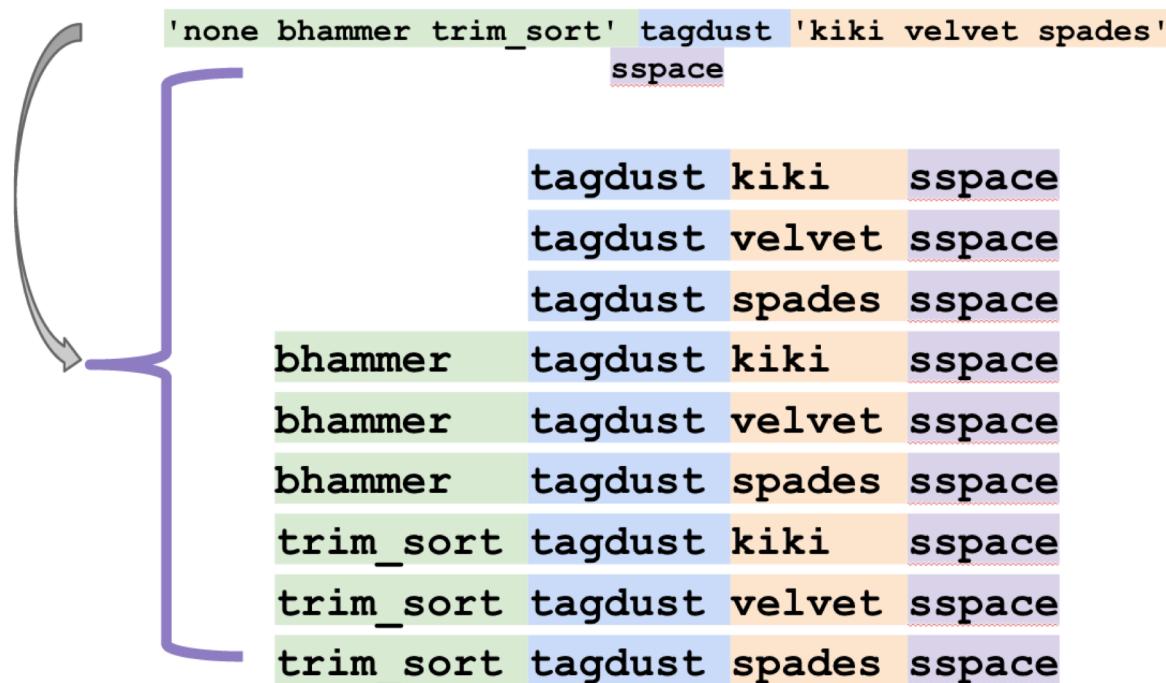
- ▶ PATRIC Team:
 - University of Chicago
 - Ryan Aydelott
 - Tom Brettin
 - Neal Conrad
 - Jim Davis
 - Emily Dietrich
 - Chris Henry
 - Dan Murphy-Olson
 - Bob Olson
 - Bruce Parrello
 - Maulik Shukla
 - Rick Stevens
 - Fangfang Xia
 - FIG
 - Terry Disz
 - Ross Overbeek
 - Gordon Pusch
 - Veronika Vonstein
 - VBI
 - Joseph Gabbard
 - Ron Kenyon
 - Dustin Machi
 - Chunhong Mao
 - Bruno Sobral
 - Rebecca Wattam
 - Andrew Warren
 - Rebecca Will
 - Harry Yoo
- Stevens Group
 - Chris Bun
 - Sebastien Boisvert

National Institute of Allergy and Infectious Diseases
Contract No. HHSN272201400027C

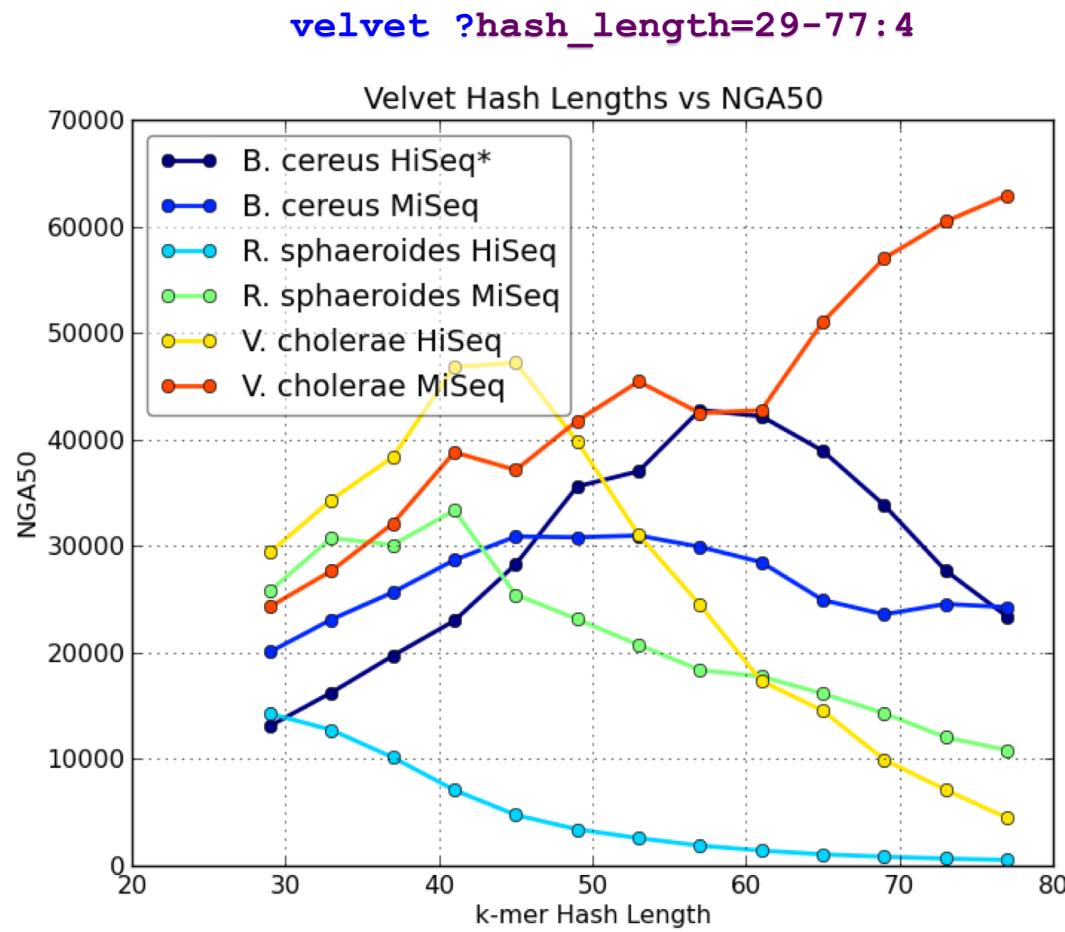
Advanced pipelines



Pipeline design



Parameter scan: k-mer Optimization



Upcoming improvements

- Improved error detection and classification using supervised learning
- Workflows for new sequencing technology (MinION), Hybrid assembly