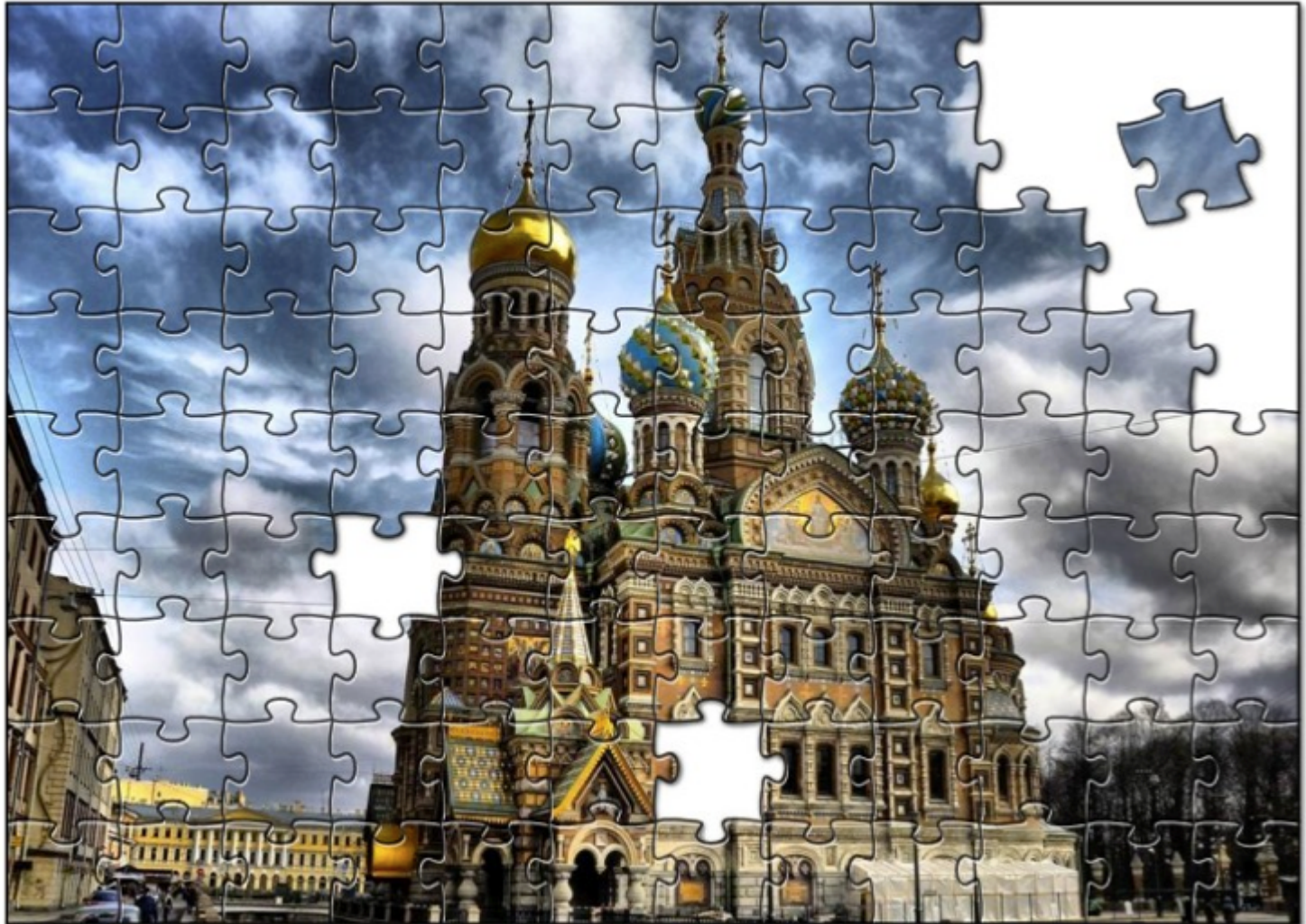


# New Frontiers of Genome Assembly with SPAdes 3.1

Andrey Prjibelski

Algorithmic Biology Laboratory, St. Petersburg Academic University,  
St. Petersburg, Russia  
<http://bioinf.spbau.ru/spades>



SPAdes (Saint Petersburg Assembler)

# SPAdes

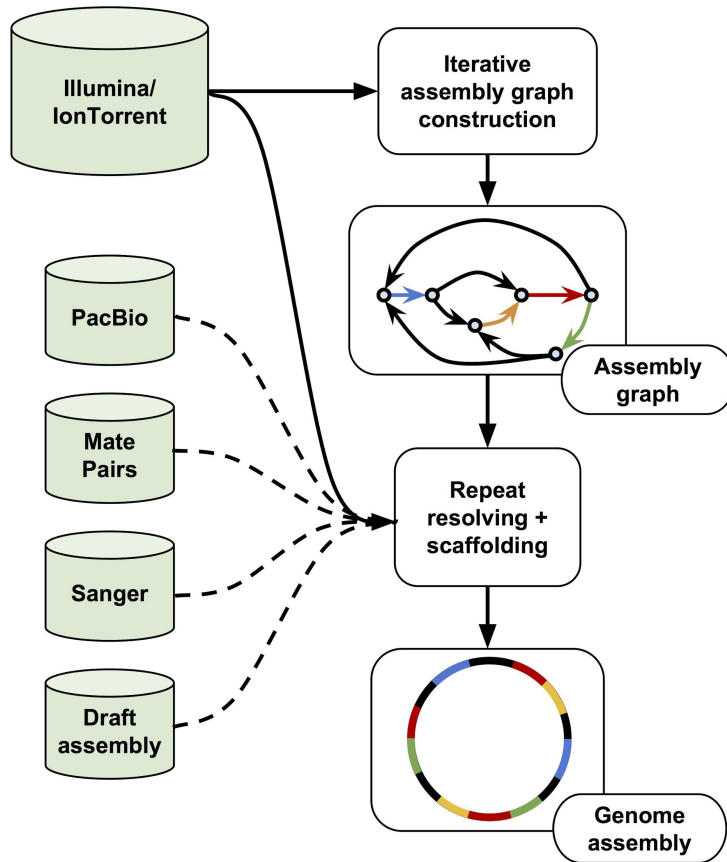
- Originally designed as single-cell assembler
- Can deal with big variations of the coverage and MDA-imposed chimeric read connections
- Turned out to work well for multi-cell isolate assemblies
- Among with MaSuRCA are top 2 assemblers in GAGE-B study by Salzberg's lab (Magoc et al, 2013)

# The Past: SPAdes 3.0

## Major additions:

- Multiple libraries support: PE, MP, single reads
- IonTorrent read error correction and assembly
- (not only) PacBio hybrid assemblies
- dipSPAdes: assembler for highly polymorphic diploid genomes

# Hybrid Assemblies



- New universal repeat resolution algorithm
- Can utilize different sources of genomic distance information

# dipSPAdes

The first de Bruijn graph assembler designed for highly polymorphic diploid genomes:



*Fungus*

heterozygosity up to 20%



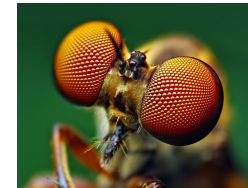
*Sea squirts*

heterozygosity up to 12%



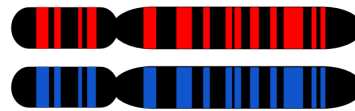
*Plants*

avg heterozygosity 7%

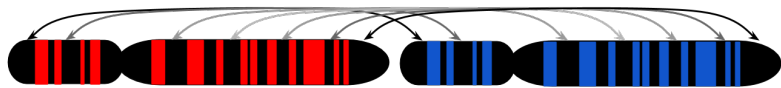


*Insects*

avg heterozygosity 9%



diploid genome with high heterozygosity



conventional approaches assemble such genome as two highly repetitive sequences and construct very fragmented assemblies



dipSPAdes takes advantage of the assembly graph structure and constructs longer consensus contigs



# IonTorrent Error Correction & Assembly

IonHammer: read error correction

- Bases on ideas of BayesHammer
- Works in homopolymer space
- Corrects both indels and mismatches
- Knows about undercall-overcall pairs



This enables accurate assemblies with inaccurate IonTorrent reads

# IonTorrent Error Correction & Assembly



E. coli O157:H7 str. Sakai from IonCommunity



5-8 fold reduction of both indel and mismatch rates!



# The Present: SPAdes 3.1

- Nextera mate pairs support and mate pair-only assemblies
- Better scaffolding and repeat resolution
- Runtime & RAM consumption improvements
- Further IonTorrent improvements
- Integration with cloud services

# The Present: SPAdes 3.1

- Nextera mate pairs support and mate pair-only assemblies
- Better scaffolding and repeat resolution
- Runtime & RAM consumption improvements
- Further IonTorrent improvements
- Integration with cloud services



Available today at <http://bioinf.spbau.ru>

# SPAdes Goes Cloud: BaseSpace

- SPAdes tuned to Illumina BaseSpace platform
- Integrates SPAdes and QUAST
- Single push-button de novo assembly solution

Check it out: <https://basespace.illumina.com/apps/160160>



# SPAdes Goes Cloud: BaseSpace



**SPADES ASSEMBLER 3.0 BETA**

Algorithmic Biology Lab

Analysis Name:

SPAdes Assembler 3.0 BETA

Sample:

Select Sample ⓘ

Save Results To:

Select a Project

▼ Assembler Options

Running Mode:

Error Correction & Assembly

Dataset Type:

Single Cell (from MDA)  
✓ Multi Cell (Isolate)

k-mer lengths:

auto ⓘ

# SPAdes for TorrentServer

- AssemblerPlus TS plugin
- Integrates SPAdes and QUAST
- Native support of BAM files as input
- Read error correction additionally uses raw flow information

# SPAdes for Cloud Platforms

SPAdes runs on:

- Illumina BaseSpace
- DNAnexus
- TorrentServer
- Galaxy (available from Galaxy Tool Shed)

# Runtime Improvements

## Improvements to BayesHammer (error correction)

- On 100 Mbp diploid genome:
  - Was: 90 hours
  - Now: 16 hours

## Improvements to repeat resolution

- On 60 Mbp repeat-rich genome:
  - Was: 78 hours
  - Now: 2 hours

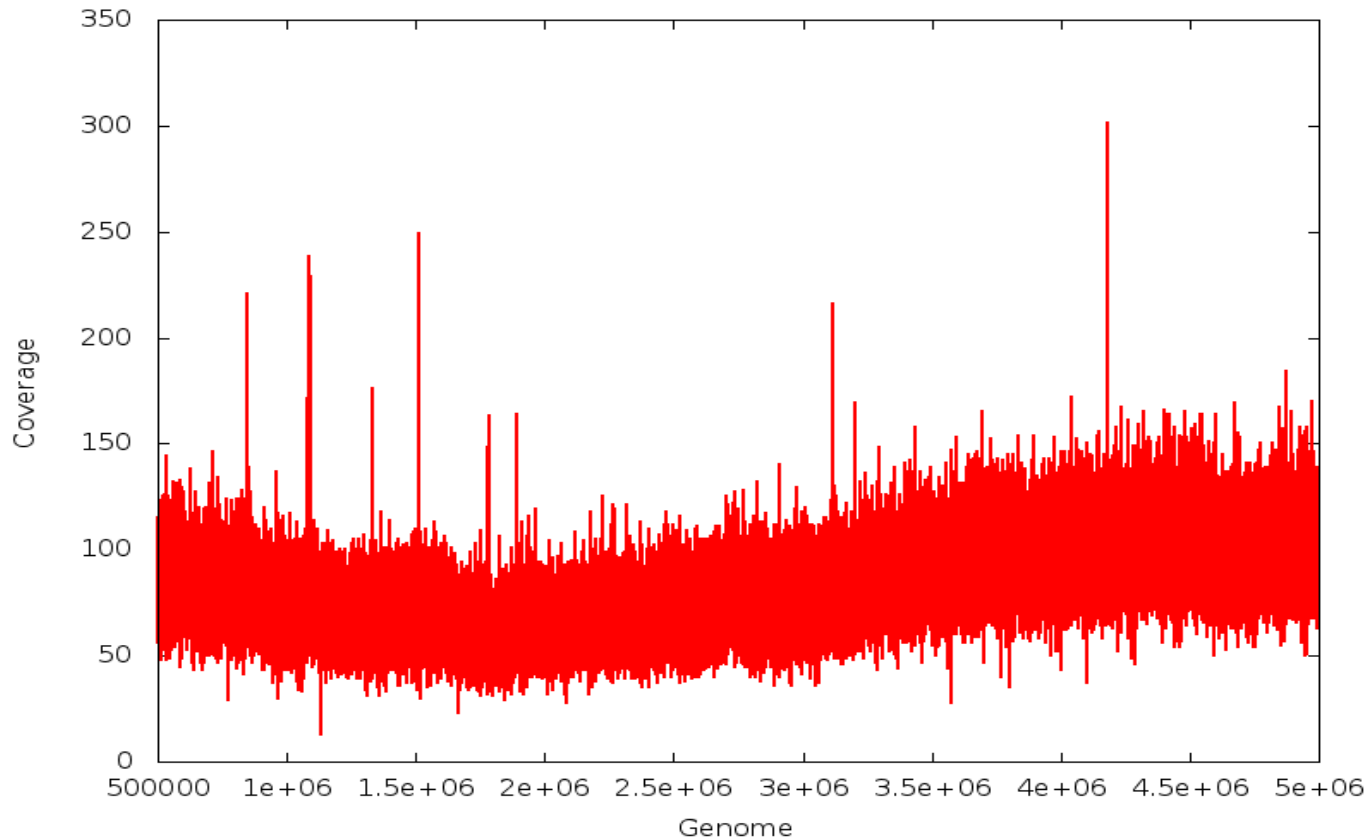
# Illumina Nextera Mate Pairs

Reasonably uniform coverage



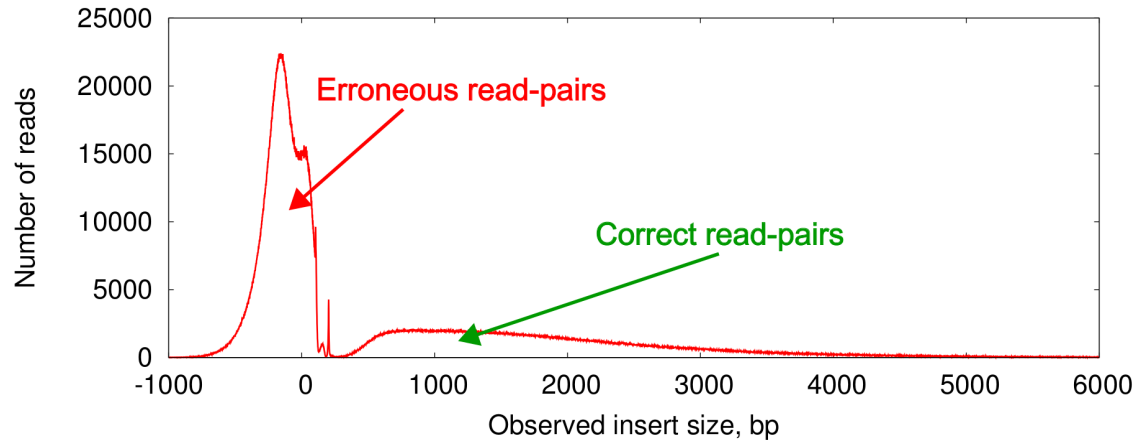
# Illumina Nextera Mate Pairs

Reasonably uniform coverage:

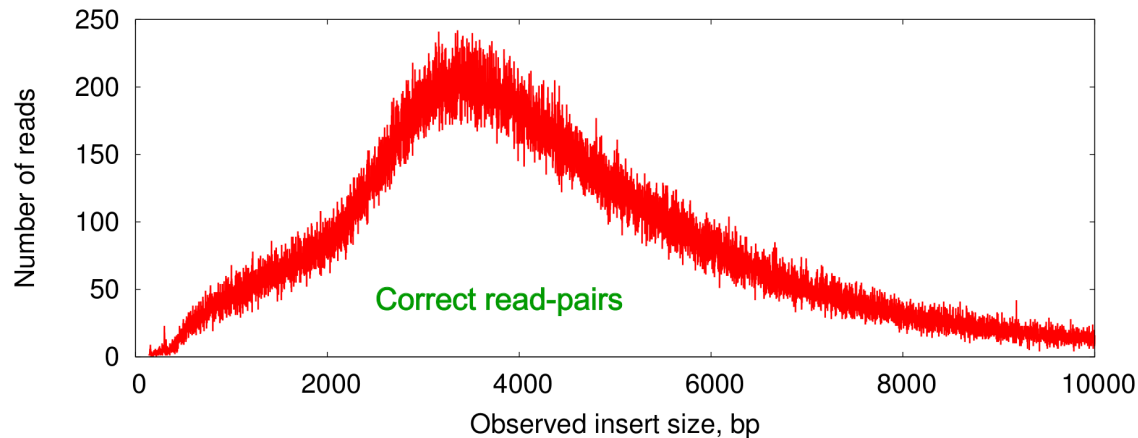


# Illumina Nextera Mate Pairs

Conventional mate-pairs:

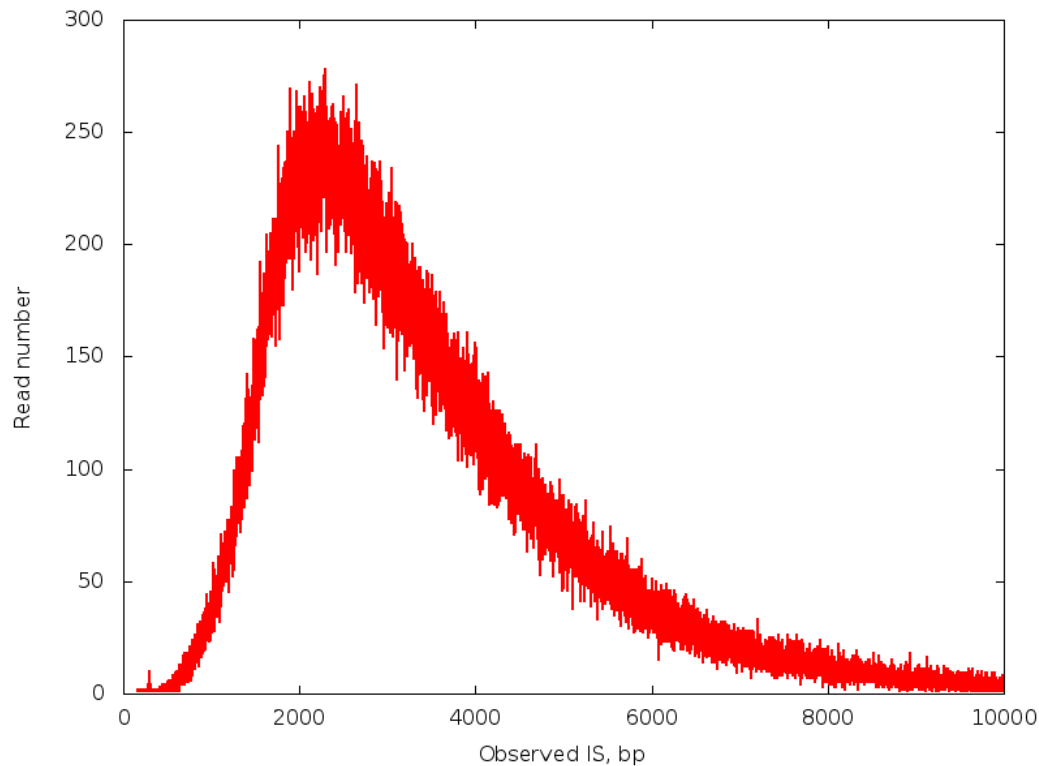


Nextera mate-pairs:



# Illumina Nextera Mate Pairs

No “paired-end” fragments:



Everything looks like large insert-size paired-end library!

# Mate Pair Only Assemblies

	SPAdes	Velvet, k=67	Velvet, k=99
# scaffolds	11	5	6
Largest contig	1496k	2817k	2006k
Reference length	2944k	2944k	2944k
N50	1496k	2817k	2006k
N75	739k	2817k	814k

*L. monocytogenes* data from Illumina



Why one would need SPAdes then?

# Mate Pair Only Assemblies

	<b>SPAdes</b>	<b>Velvet, k=67</b>	<b>Velvet, k=99</b>
# scaffolds	11	5	6
Reference length	2944k	2944k	2944k
# misassemblies	<b>0</b>	<b>6</b>	<b>3</b>
Largest alignment	<b>1496k</b>	1493k	1493k
NGA50	<b>1496k</b>	1493k	1493k
NGA75	<b>739k</b>	505k	509k
Genome fraction	<b>99.463%</b>	99.133%	99.391%

# Acknowledgement

## SPAdes team:

Dmitry Antipov

Anton Bankevich

Sergey Nurk

Alexey Gurevich

Anton Korobeynikov

Yana Safonova

Irina Vasilinetc

Alla Lapidus

Pavel Pevzner



This work was supported by the Government of the Russian Federation (grant 11.G34.31.0018)

# Thank you!

SPAdes: [bioinf.spbau.ru/spades](http://bioinf.spbau.ru/spades)

QUAST: [bioinf.spbau.ru/quast](http://bioinf.spbau.ru/quast)

