PALMapper: Fast, Accurate and Variation-aware RNA-seq Alignments

<u>André Kahles,</u> Géraldine Jean, David Kuo, Vipin T. Sreedharan, Gunnar Rätsch



New York City, USA

HiTSeq, Berlin, July 20, 2013















Applications for Variation-aware Alignments



Resource Name	Organism	Strains/Individuals	SNPs/SNVs	InDels	Reference
Mouse Genome Project	M. musculus	17 strains	129,260,574	21,683,297	Keane (2011)
19 Genomes Project	A. thaliana	19 strains	3,070,000	1,200,000	Gan (2011)
Panzea	T. mays	103 strains	55,061,920	3,200,000	Chia (2012)
Million Mutation Project	C. elegans	2,047 strains	800,000	220,000	Thompson (2013)
1000 Genomes Project	H. sapians	1,092 individuals	36,700,000	1,380,000	The 1000 Genomes
					Project Consortium (2012)

General Workflow

- general seed-and-extend strategy
- following four basic steps

General Workflow

- general seed-and-extend strategy
- following four basic steps



Index and Seeding (normal case)

- Indexing genome with k-mer based index or bwt-based index
- Lookup read *k*-mers in the genome index
- Generate long and short seed regions for further alignment.



HTS Reads



Index and Seeding (normal case)

- Indexing genome with k-mer based index or bwt-based index
- Lookup read *k*-mers in the genome index
- Generate long and short seed regions for further alignment.



HTS Reads



- Indexing genome with k-mer based index or bwt-based index
- Lookup read *k*-mers in the genome index
- Generate long and short seed regions for further alignment



- Indexing genome with k-mer based index or bwt-based index
- Lookup read *k*-mers in the genome index
- Generate long and short seed regions for further alignment



- Indexing genome with k-mer based index or bwt-based index
- Lookup read k-mers in the genome index
- Generate long and short seed regions for further alignment



Index and Seeding (normal case)

- Indexing genome with k-mer based index or bwt-based index
- Lookup read k-mers in the genome index
- Generate long and short seed regions for further alignment



ACCGTCGCGCGCGCG...AGAACGCT

- Indexing genome with k-mer based index or bwt-based index
- Lookup read k-mers in the genome index
- Generate long and short seed regions for further alignment



- Indexing genome with k-mer based index or bwt-based index
- Lookup read *k*-mers in the genome index
- Generate long and short seed regions for further alignment



Indexing (variant case)

 Include available variation into index (dbSNP, COSMIC, ENCODE, 1000 Genomes . . .

• Principle for querying the index stays the same



- ACCGTCGCGCGCGT...TCGGCG...AGAACGCT ACCGTCGAGCGCGT...TCGGCG...AGAACGCT

HTS Reads



Indexing (variant case)

- Include available variation into index (dbSNP, COSMIC, ENCODE, 1000 Genomes ...)
- Principle for querying the index stays the same



Indexing (variant case)

- Include available variation into index (dbSNP, COSMIC, ENCODE, 1000 Genomes ...)
- Principle for querying the index stays the same



Indexing (variant case)

- Include available variation into index (dbSNP, COSMIC, ENCODE, 1000 Genomes ...)
- Principle for querying the index stays the same



Additional to SNPs PALMapper can account for

- Insertions
- Deletions

• Combinations of variants (incl. unbalanced substitutions)

Insertions

ACCG---TCGCGAGCGT...TCGGCG...AGAACGCT ACCGCGTTCGCGAGCGT...TCGGCG...AGAACGCT

G---TCG GCGT CGTT GTTC TTCG

Additional to SNPs PALMapper can account for

- Insertions
- Deletions

Combinations of variants (incl. unbalanced substitutions)



Additional to SNPs PALMapper can account for

- Insertions
- Deletions
- Combinations of variants (incl. unbalanced substitutions)



Projection to Pseudo Chromosomes

- Rank and combine seed regions from previous step
- Integrate with splice site prediction and coverage information
- Project to genome and form Pseudo Chromosome Sequences



Projection to Pseudo Chromosomes

- Rank and combine seed regions from previous step
- Integrate with splice site prediction and coverage information
- Project to genome and form Pseudo Chromosome Sequences



Projection to Pseudo Chromosomes

- Rank and combine seed regions from previous step
- Integrate with splice site prediction and coverage information
- Project to genome and form Pseudo Chromosome Sequences



Projection to Pseudo Chromosomes

- Rank and combine seed regions from previous step
- Integrate with splice site prediction and coverage information
- Project to genome and form *Pseudo Chromosome Sequences*



Projection to Pseudo Chromosomes

- Rank and combine seed regions from previous step
- Integrate with splice site prediction and coverage information
- Project to genome and form *Pseudo Chromosome Sequences*



Projection to Pseudo Chromosomes (variant case)

Include known sequence variants to form combination of pseudo chromosomes



Dynamic Program Alignment Against Graph

- Integrate substitutions as IUPAC codes for ambiguous bases
- Integrate deletions as jumps
- Integrate insertions as jumps
- Collapse all variations into graph-like structure
- Align against variation graph

Genome 🗔

ACCGTACGGT

Dynamic Program Alignment Against Graph

Integrate substitutions as IUPAC codes for ambiguous bases

- Integrate deletions as jumps
- Integrate insertions as jumps
- Collapse all variations into graph-like structure
- Align against variation graph

~	
Genome	

Substitutions

ACCGTACGGT AC**G**GTAC**A**GT

Dynamic Program Alignment Against Graph

- Integrate substitutions as IUPAC codes for ambiguous bases
- Integrate deletions as jumps
- Integrate insertions as jumps
- Collapse all variations into graph-like structure
- Align against variation graph

Genome	ACCGTACGGT
ubstitutions	AC G GTAC A GT
Deletions	ACCGCGGT

S

Dynamic Program Alignment Against Graph

- Integrate substitutions as IUPAC codes for ambiguous bases
- Integrate deletions as jumps
- Integrate insertions as jumps
- Collapse all variations into graph-like structure
- Align against variation graph



Dynamic Program Alignment Against Graph

- Integrate substitutions as IUPAC codes for ambiguous bases
- Integrate deletions as jumps
- Integrate insertions as jumps
- Collapse all variations into graph-like structure

• Align against variation graph



Dynamic Program Alignment Against Graph

- Integrate substitutions as IUPAC codes for ambiguous bases
- Integrate deletions as jumps
- Integrate insertions as jumps
- Collapse all variations into graph-like structure
- Align against variation graph





Dynamic Program (DP)

The inferred seed region guides a DP-based alignment algorithm:



Semi-Global	Banded	Banded Spliced	
Align whole read	Limit # of gaps	Allow long gaps	Account for
to part of pseudo	from perfect	corresponding to	genomic
chromosome	alignment	introns	variations
André Kahles (SKI, New York)	PALMapper		liTSea, July 20, 2013 10

Variant alignments cause gain in sensitivity

• RNA-Seq of A. thaliana strain Can-0 aligned to reference Col-0

use variation during seeding and extension improves sensitivity



Variant alignments cause gain in sensitivity

• RNA-Seq of A. thaliana strain Can-0 aligned to reference Col-0

use variation during seeding and extension improves sensitivity



- RNA-Seq of A. thaliana strain Can-0 aligned to reference Col-0
- use variation during seeding and extension improves sensitivity



- RNA-Seq of A. thaliana strain Can-0 aligned to reference Col-0
- use variation during seeding and extension improves sensitivity



- RNA-Seq of A. thaliana strain Can-0 aligned to reference Col-0
- use variation during seeding and extension improves sensitivity



- RNA-Seq of A. thaliana strain Can-0 aligned to reference Col-0
- use variation during seeding and extension improves sensitivity



Running Time

Running time and Scaling

- 2M reads / h / core on human genome
- ca 1M reads / h /core with variants (complexity dependent)
- almost linear scaling with number of cores

Accuracy of Intron Prediction

Comparison of predicted introns to annotation

• PALMapper outperforms TopHat and STAR in terms of accuracy

Variation-aware alignment increases performance



Accuracy of Intron Prediction

Comparison of predicted introns to annotation

- PALMapper outperforms TopHat and STAR in terms of accuracy
- Variation-aware alignment increases performance



- Use FluxSimulator to generate artificial reads (10M from 5K genes)
- Duplicate and mutate read set
- Re-align to original genome
- Expectation: perfect allelic balance!



- Use FluxSimulator to generate artificial reads (10M from 5K genes)
- Duplicate and mutate read set
- Re-align to original genome
- Expectation: perfect allelic balance!



- Use FluxSimulator to generate artificial reads (10M from 5K genes)
- Duplicate and mutate read set
- Re-align to original genome
- Expectation: perfect allelic balance!



- Use FluxSimulator to generate artificial reads (10M from 5K genes)
- Duplicate and mutate read set
- Re-align to original genome
- Expectation: perfect allelic balance!



Allelic Imbalance

Recovery of perfect allelic balance in artificial set

heterozygous alleles are expected to be balanced



Measure differential gene expression with DESeq



Measure differential gene expression with DESeq



Measure differential gene expression with DESeq



Measure differential gene expression with DESeq



Versatile Set of Additional Features

- Fully parallelized (used tested on 64 core machines)
- Various read trimming: 3' end and polyA-tail
- Handles strand-specific reads
- Built-in filtering
- Allows to report arbitrary many multi-mappers
- Supports non-canonical splice sites
- Can align over several introns
- Junction remapping
- Multimapper resolution
- ..

Versatile Set of Additional Features

- Fully parallelized (used tested on 64 core machines)
- Various read trimming: 3' end and polyA-tail
- Handles strand-specific reads
- Built-in filtering
- Allows to report arbitrary many multi-mappers
- Supports non-canonical splice sites
- Can align over several introns
- Junction remapping
- Multimapper resolution
- ..



Availability

Command-line

http://github.org/ratschlab/palmapper
http://ftp.raetschlab.org/software/palmapper/

- Galaxy web-interface http://galaxy.cbio.mskcc.org/
- Tutorial paper: G. Jean et al., Current Protocols in Bioinformatics, 2010.
- Original paper: A. Kahles & G. Jean et al., forthcoming.

Acknowledgements





Gunnar Rätsch Supervision

Geraldine Jean Variation-aware DP

David Kuo Variant Parsing



Vipin T. Sreedharan Galaxy Integration

Funding by DFG & MSKCC.

Thank you for your attention.