

# Mapping of bisulfite sequencing data

Christian Otto

Bioinformatics, Leipzig

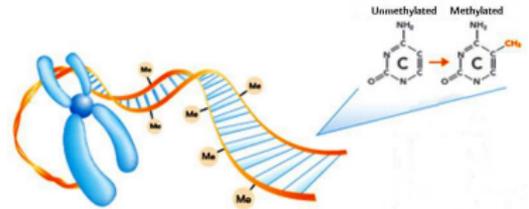
February 2011

# Table of content

- 1 Biological background
- 2 Mapping
- 3 Results
- 4 Future work

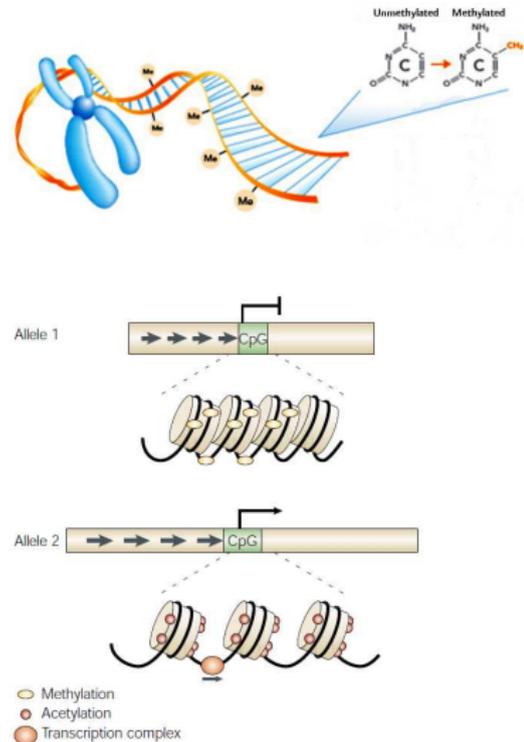
# DNA methylation

- methylation on C5 position of cytosines
- occurs at CpG dinucleotides and at non-CpG dinucleotides in plants and embryonic stem cells in Human



# DNA methylation

- methylation on C5 position of cytosines
- occurs at CpG dinucleotides and at non-CpG dinucleotides in plants and embryonic stem cells in Human
- hypermethylation of promoters is correlated with heterochromatin formation and silenced transcription
- vital role in:
  - embryonic development
  - maintenance of pluripotency
  - X-chromosome inactivation
  - genomic imprinting

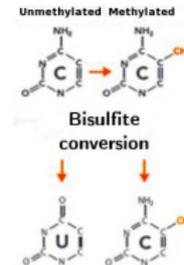


# Techniques to capture DNA methylation

- 1 Methylated DNA immunoprecipitation sequencing (MeDIP-seq):
  - anti-methylcytosine antibody
  - sequencing of captured fragments
  - common resolution of 150bp

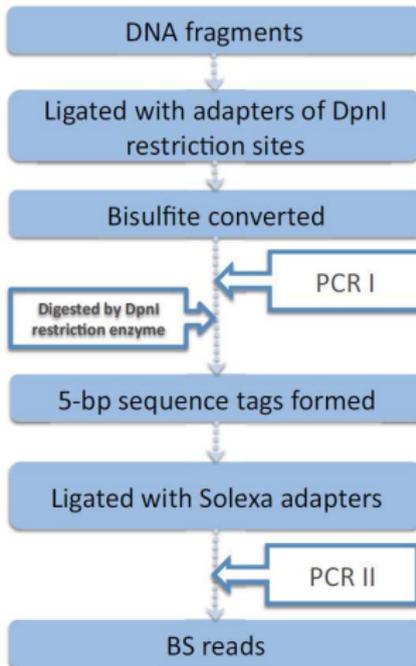
# Techniques to capture DNA methylation

- 1 Methylated DNA immunoprecipitation sequencing (MeDIP-seq):
  - anti-methylcytosine antibody
  - sequencing of captured fragments
  - common resolution of 150bp
- 2 Bisulfite sequencing:
  - treatment with sodium bisulfite
  - conversion of unmethylated Cs to Us
  - sequencing of converted fragments
  - **adjustments in mapping algorithms**
  - single-base resolution  $\Rightarrow$  **gold standard**
  - drawback: very costly

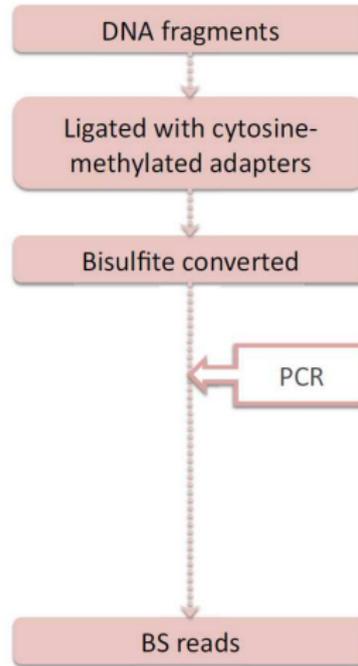


# Bisulfite sequencing

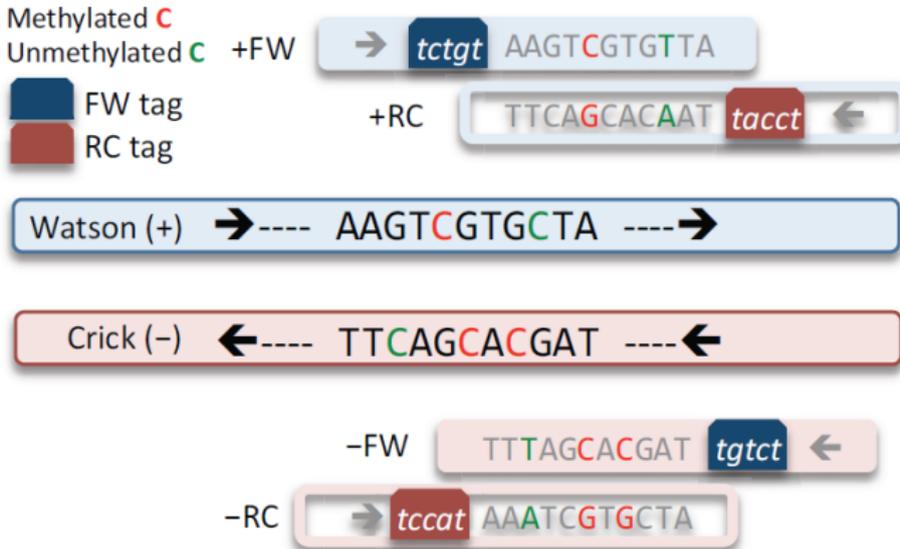
## Cokus *et al*'s library protocol



## Lister *et al*'s library protocol



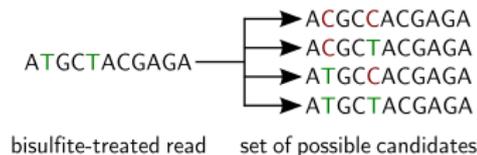
# Bisulfite sequencing (cont'd)



⇒ needs to be taken into account by mapping algorithms

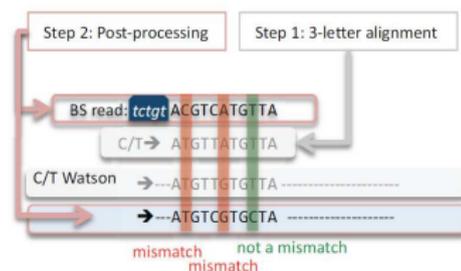
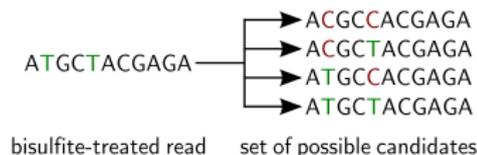
# Previous approaches

- 1 enumeration:
  - generate set of possible candidates that can result in the bisulfite read
  - very time-consuming
  - e.g. BSMAP



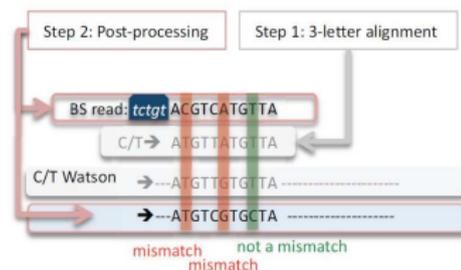
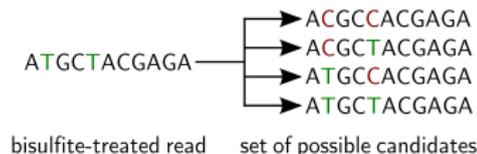
# Previous approaches

- 1 enumeration:
  - generate set of possible candidates that can result in the bisulfite read
  - very time-consuming
  - e.g. BSMAP
- 2 collapsing alphabet:
  - convert Cs in reference and query sequence to Ts
  - no distinction between  $C \rightarrow T$  and  $T \rightarrow C$  mismatches possible
  - e.g. BS Seeker, MAQ



# Previous approaches

- 1 enumeration:
  - generate set of possible candidates that can result in the bisulfite read
  - very time-consuming
  - e.g. BSMAP
- 2 collapsing alphabet:
  - convert Cs in reference and query sequence to Ts
  - no distinction between C→T and T→C mismatches possible
  - e.g. BS Seeker, MAQ
- 3 wildcard matching:
  - allow only bisulfite mismatches
  - e.g. RMAP



# segemehl's approach

bisulfite-treated read    GGATATCGATGTTACGAGTTCGTTT

methyated sites

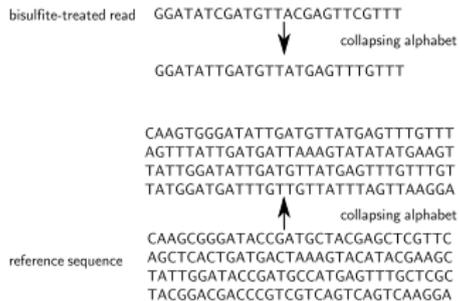
The diagram shows a bisulfite-treated read sequence: GGATATCGATGTTACGAGTTCGTTT. The 'C' characters at positions 5, 10, and 15 are highlighted with red lines. Three arrows point from the text 'methyated sites' above to these three 'C' characters.

reference sequence

```
CAAGCGGGATACCGATGCTACGAGCTCGTTC
AGTCTACTGATGACTAAAGTACATACGAAGC
TATTGGATACCGATGCCATGAGTTTGCTCGC
TACGGACGCCGTCTCAGTCAGTCAAGGA
```

Hybrid approach:

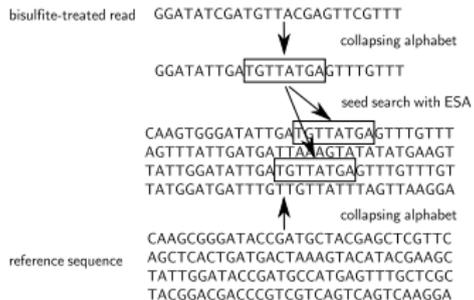
# segemehl's approach



Hybrid approach:

- 1 seed search in ESA on collapsed alphabet

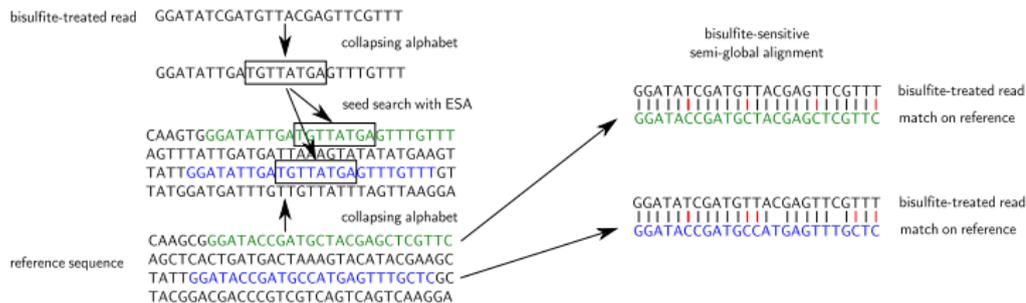
# segemehl's approach



Hybrid approach:

- 1 seed search in ESA on collapsed alphabet

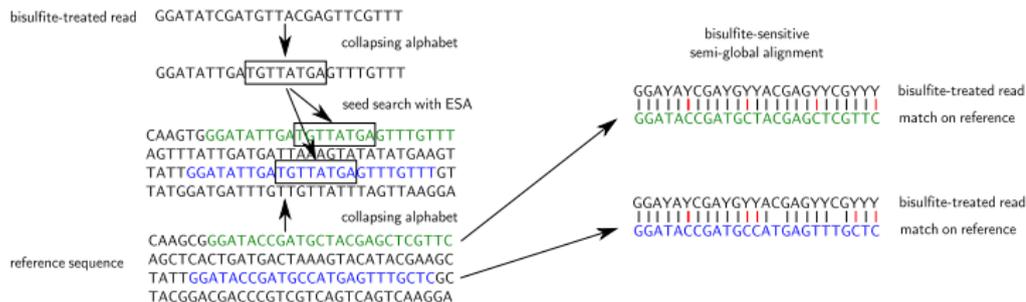
# segemehl's approach



## Hybrid approach:

- 1 seed search in ESA on collapsed alphabet
- 2 semi-global alignment allowing only bisulfite mismatches using Myers bitvector algorithm

# segemehl's approach



## Hybrid approach:

- 1 seed search in ESA on collapsed alphabet
- 2 semi-global alignment allowing only bisulfite mismatches using Myers bitvector algorithm

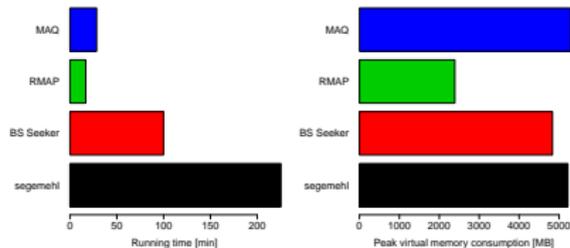
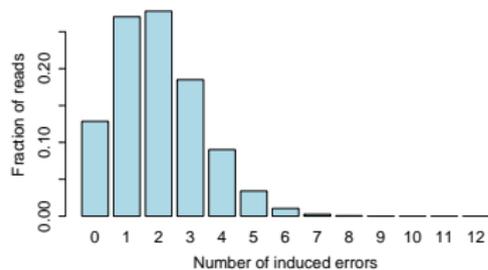
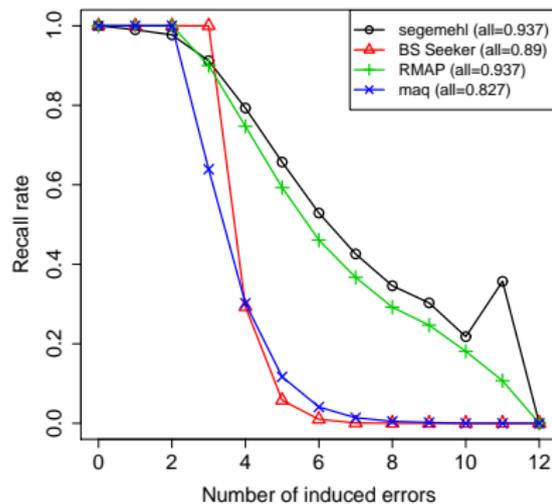


# Artificial datasets

- 200MB random reference
- 10 million reads (40nt or 80nt)
- methylC-seq protocol with 50% methylation rate
- induced errors (5% or 10% error rate)

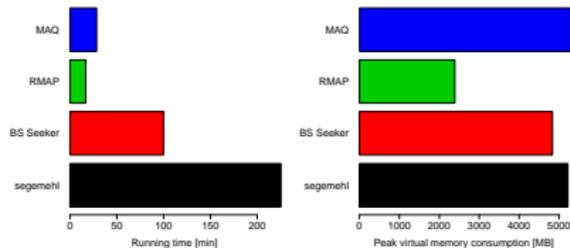
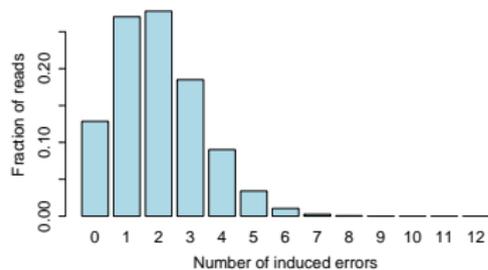
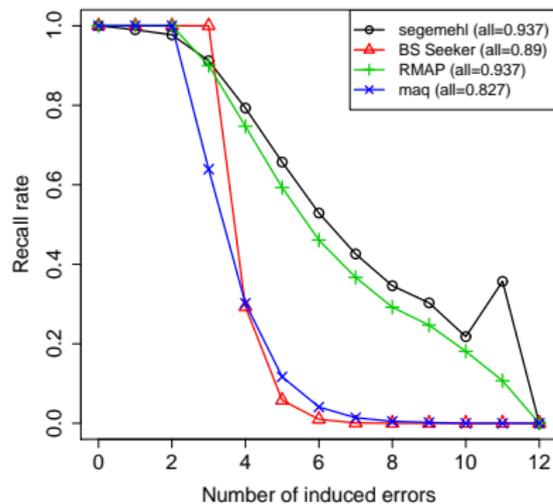
# Performance on artificial data

(1) short reads, 5% mismatches



# Performance on artificial data

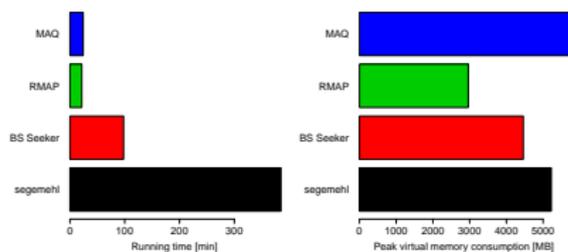
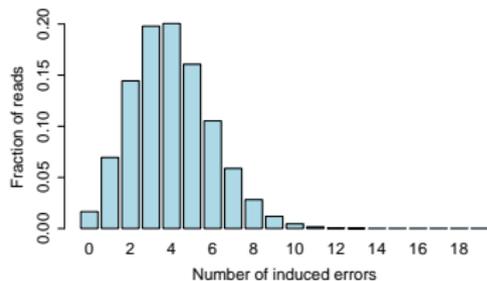
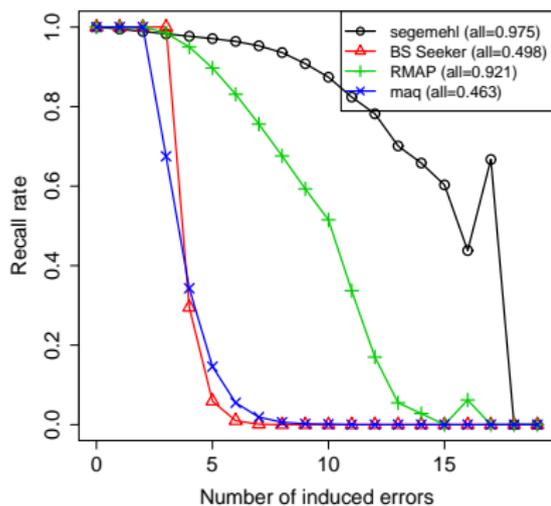
(1) short reads, 5% mismatches



⇒ similar recall with all tools

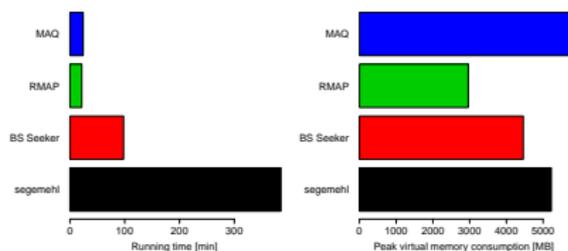
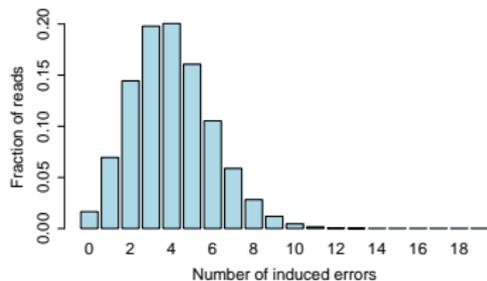
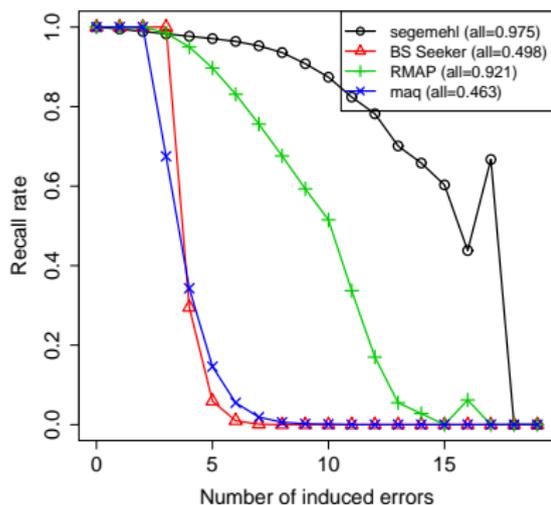
# Performance on artificial data

(1) longer reads, 5% mismatches



# Performance on artificial data

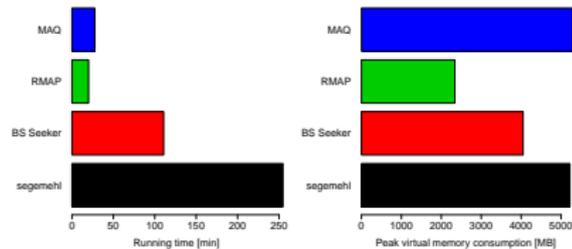
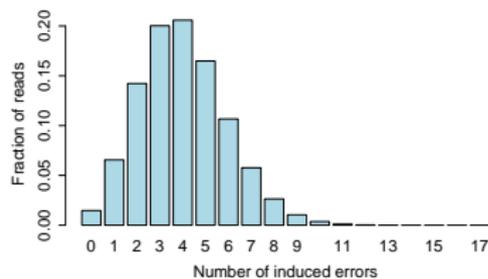
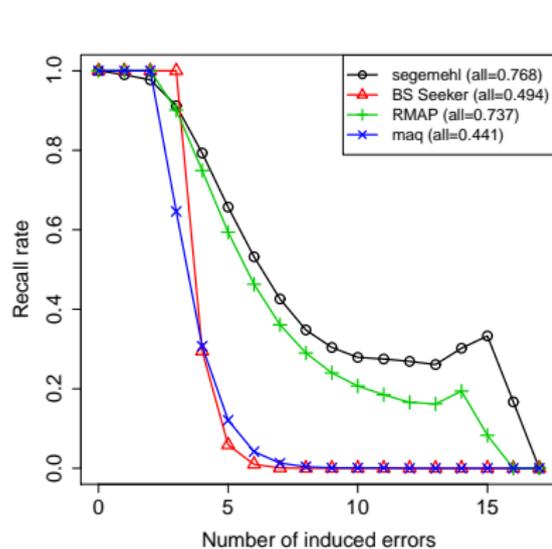
(1) longer reads, 5% mismatches



⇒ gain in recall with `segemehl` but decline with other tools

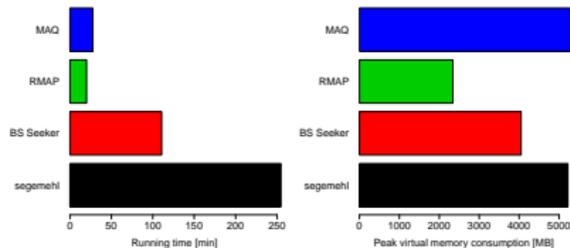
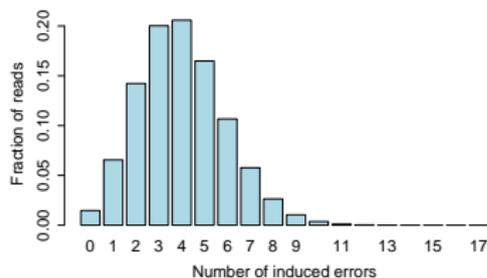
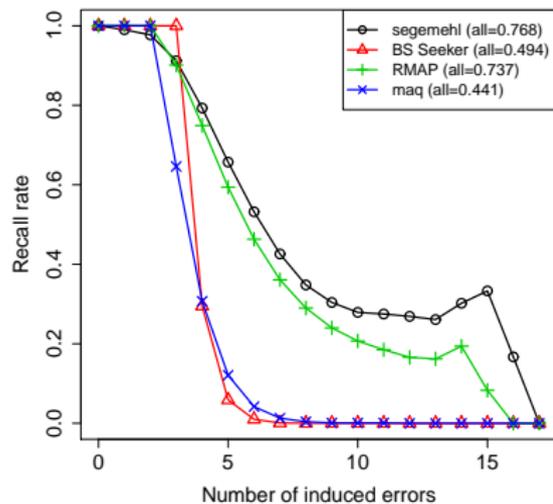
# Performance on artificial data

(1) short reads, 10% mismatches



# Performance on artificial data

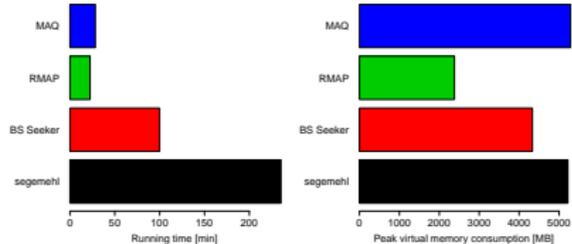
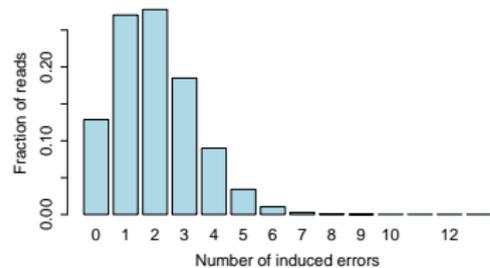
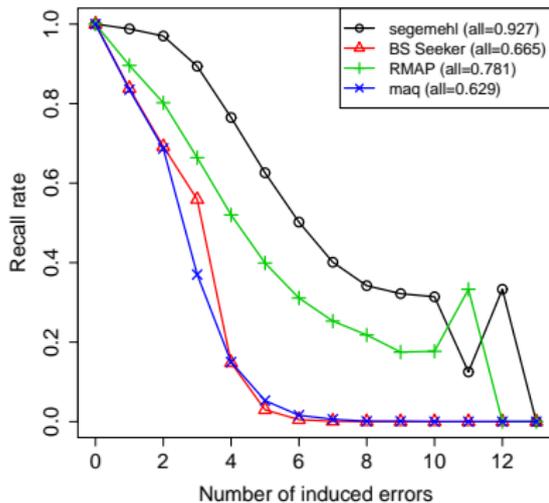
(1) short reads, 10% mismatches



⇒ all tools achieve lower recall

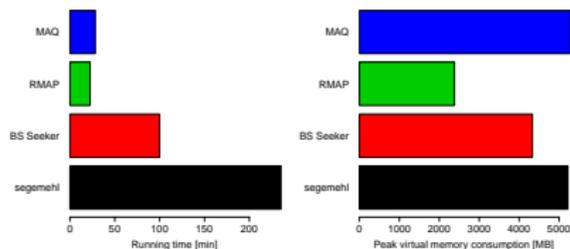
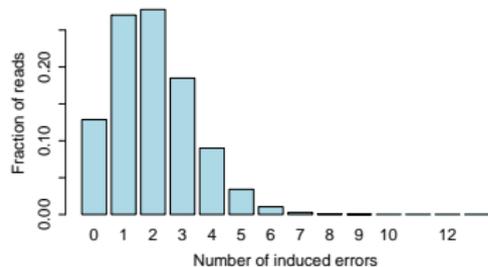
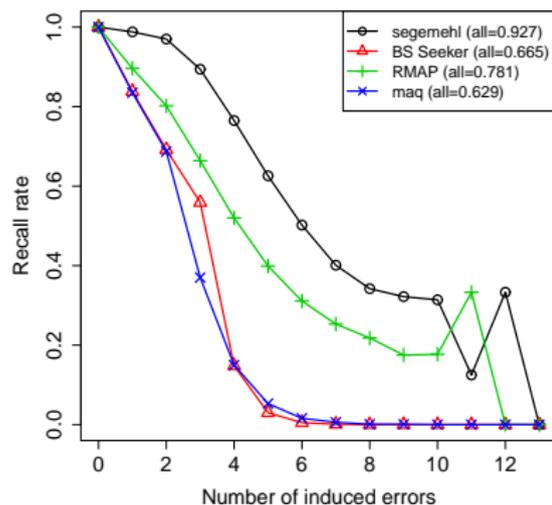
# Performance on artificial data

(1) short reads, 5% errors (mismatches, insertions, deletions)



# Performance on artificial data

(1) short reads, 5% errors (mismatches, insertions, deletions)



⇒ only segemehl takes indels into account

# Real-life dataset

SRR019048<sup>1</sup>:

- whole genome shotgun bisulfite sequencing of the Human H1 cell line
- $\approx$  15 million reads of length 87nt
- methylC-seq protocol used

---

<sup>1</sup>Lister et al. **Human DNA methylomes at base resolution show widespread epigenomic differences.** Nature (2009)

# Real-life dataset

SRR019048<sup>1</sup>:

- whole genome shotgun bisulfite sequencing of the Human H1 cell line
- $\approx$  15 million reads of length 87nt
- methylC-seq protocol used

Mapping results:

program	uniquely/best mapped reads	running time (in min)	peak virtual memory (in MB)
segemehl	13'367'984 (87.2%)	816	74453.60
BS Seeker	6'243'531 (40.7%)	247	9280.45
RMAP	9'243'240 (60.3%)	962	7716.59
MAQ	8'723'244 (56.9%)	4731	8327.98

---

<sup>1</sup>Lister et al. **Human DNA methylomes at base resolution show widespread epigenomic differences.** Nature (2009)

# Summary

segemeh1 performs equally to other tools in case of short reads with only few uniformly distributed mismatches

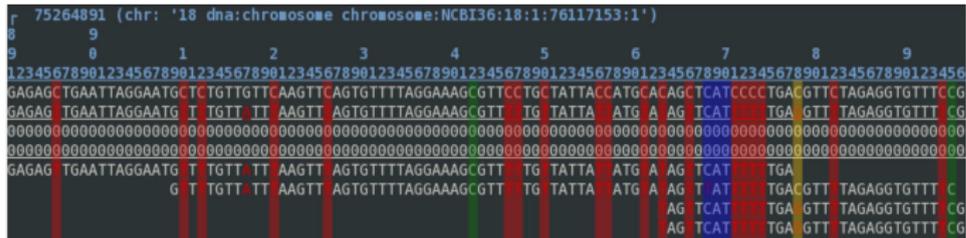
but:

- gains higher recall with longer reads
- can cope with higher error rates
- considers insertions and deletions
- supports multi-threading
- **outperforms other tools in real-life datasets**



# Postprocessing (cont'd)

Identification of methylation-relevant sites:



non-methylated site

methylated site  
(CpG)

methylated site  
(CHH)

methylated site  
(CpG) and SNP

# Postprocessing (cont'd)

Calling methylation state:

- adapt SNP calling
- assess confidence of calls
- incorporate non-uniquely mapped reads
- consider genomic surrounding
- detect monoallelic modifications

⇒ report detailed information on cytosines with highly confident methylation calls

# The end

Thank you for listening!

Feel free to ask some questions.