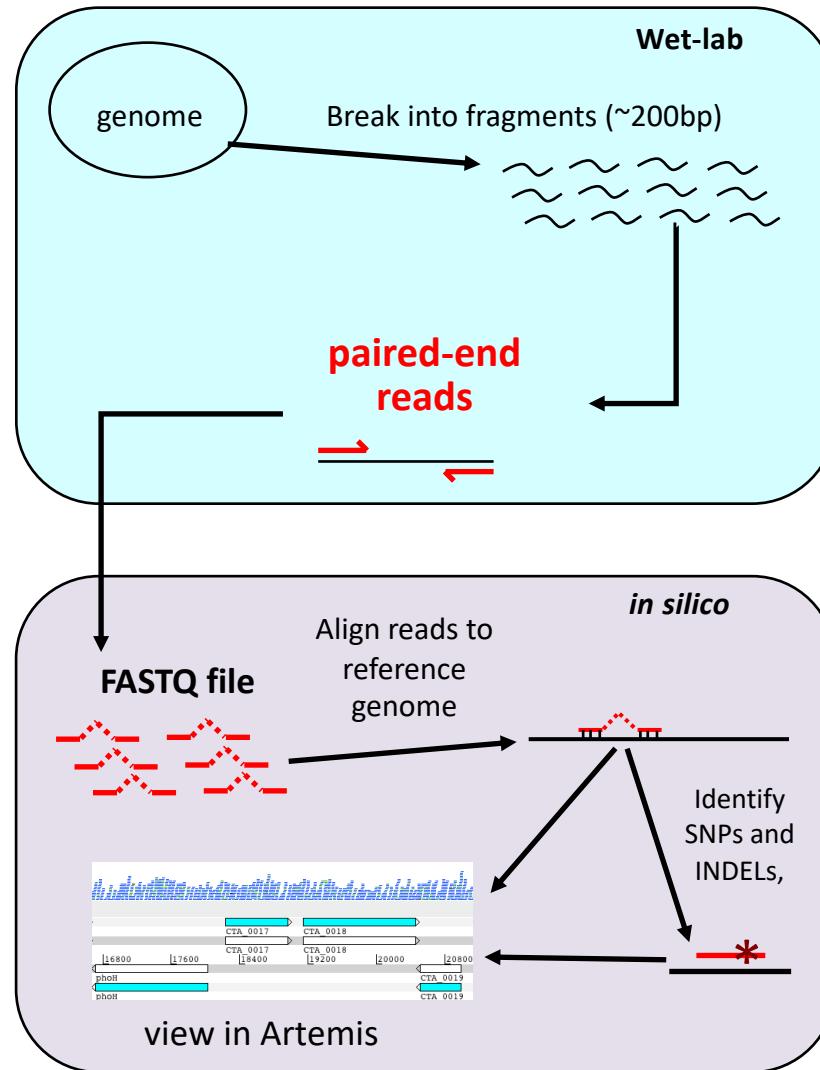


Module 2

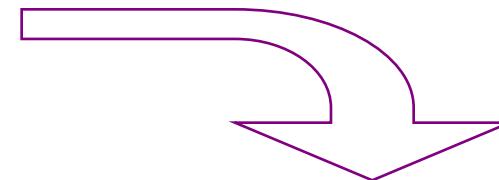
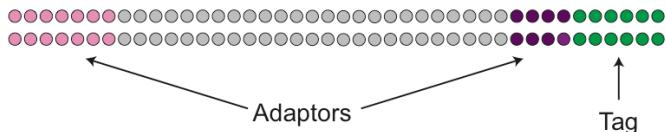
Mapping short reads

					
Key Methods	Amplicon, targeted RNA, small RNA, and targeted gene panel sequencing.	Small genome, amplicon, and targeted gene panel sequencing.	Everyday exome, transcriptome, and targeted resequencing.	Production-scale genome, exome, transcriptome sequencing, and more.	Population- and production-scale whole-genome sequencing.
Maximum Output	7.5 Gb	15 Gb	120 Gb	1500 Gb	1800 Gb
Maximum Reads per Run	25 million	25 million [†]	400 million	5 billion	6 billion
Maximum Read Length	2 × 150 bp	2 × 300 bp	2 × 150 bp	2 × 150 bp	2 × 150 bp
Run Time	4–24 hours	4–55 hours	12–30 hours	<1–3.5 days (HiSeq 3000/HiSeq 4000) 7 hours–6 days (HiSeq 2500)	<3 days

Workflow of re-sequencing, alignment, and *in silico* analysis

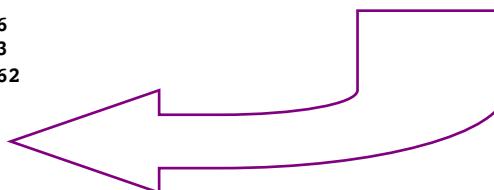
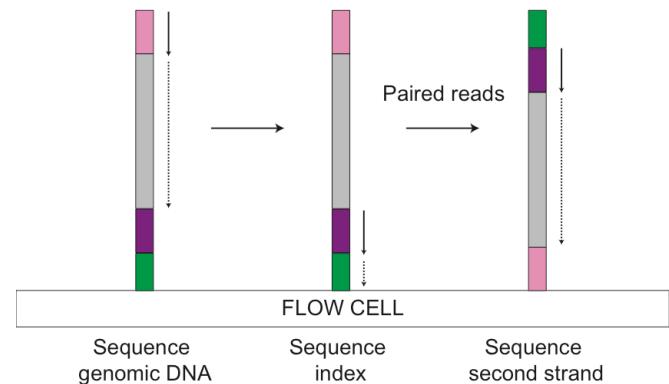


Illumina Sequencing Reads



FASTQ file

TAG	Strain
CGTGAT	HGS A942
ACATCG	3HK
GCCTAA	CHI59
TGGTCA	CHI61
CACTGT	BK2491
ATTGGC	TUR9
GATCTG	HU25
TCAAGT	HU109
CTGATC	HSJ216
AAGCTA	FFP103
GTAGCC	ICP5062
TACAAG	GRE4



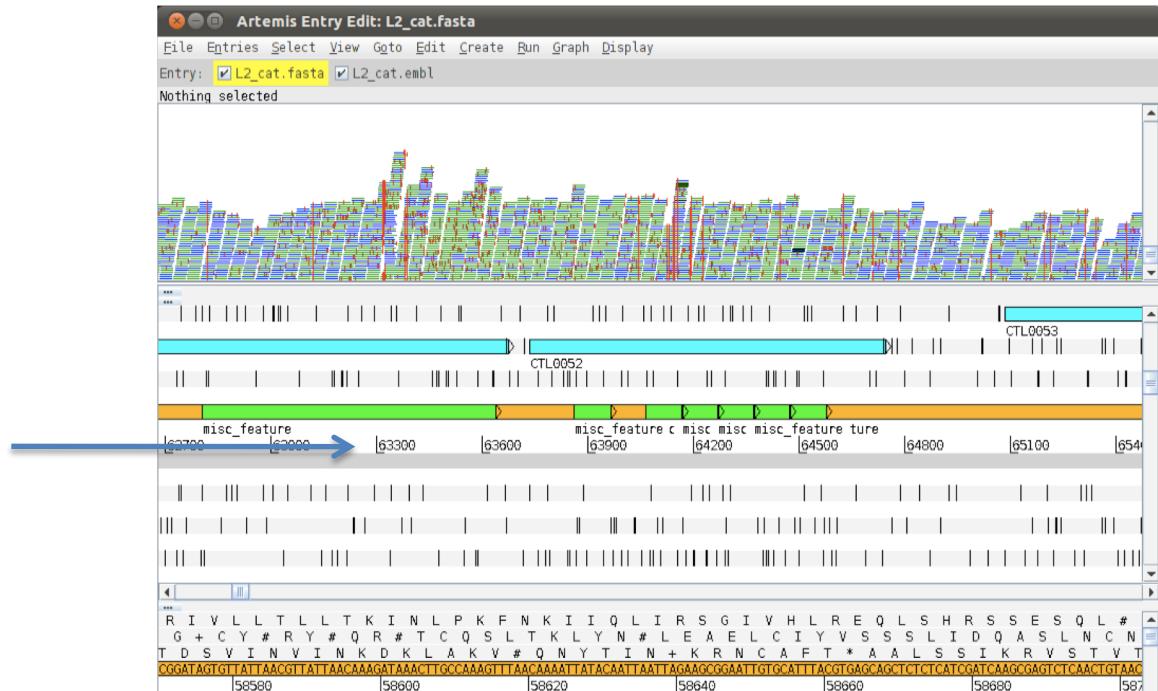
Mapping Illumina sequence data

Resequencing – takes a set of short reads and maps these onto a reference genome

```
@IL24_5151:3:1:1553:916#9/1
NAAACTACTACACCCACTCAGGACACCAGGGACATCATTGCTGACGCCACGGCCTCACAGTGCTGAGCTGATGAT
+
$705291596=>>>=>=>=>=>=>535:6=>>>=>=5:;318656:==991/1,-0,0015204.1
@IL24_5151:3:1:2173:904#9/1
NTTTAACCGTACTTCACCAGGATTATCGCAGGCGGATTCTGGTATTCAAAAATAGCGTTAATCCA
+
$948883999>=>>>=>=>9>>>=>=>=>:::==>=55:88=>9:0;;:==>=>>>5
@IL24_5151:3:1:2948:912#9/1
NCCACCAGACACTGTCCGCAACCCCGTAAGGGGCCAACGTTAGAACATCAAACATTAAAGGGTGGTATTCAAGG
+
```

>reference sequence

```
ATTGAACGCTGGCGGCAGGCCTAACACATGCAAGTC
GAGCGGCAGCGGGAAAGTAGTTT
TACTTGCCGGCGAGCGGGGACGGGTGAGTAATGTC
CTGGGAAACTGCCTGATGGAGG
GATAACTACTGAAACGGTAGCTAACCGCATGACC
TCGTAAGAGCAAAGTGGGGAC
TTCGGGCCTCACGCCATGGATGTGCCAGATGGGAT
TAGCTAGTAGGTGGGTAATGG
TCACCTAGGCAGCATCCCTAGCTGGTCTGAGAGGAT
GACCAGCCACACTGGAACGTGAG
CACGGTCCAGACTCCTACGGGAGGCAGCAGTGGGA
ATATTGCACAATGGCGCAAGC
GATGCAGCCATGCCGCGTGTGAAGAAGGCCTCG
GGTTGAAAGCACTTCAGCGAG
AGGAAGGCAGTCGTGTTAATGACGATTGATTGAC
GTTACTCGCAGAAGAACCGGGC
```

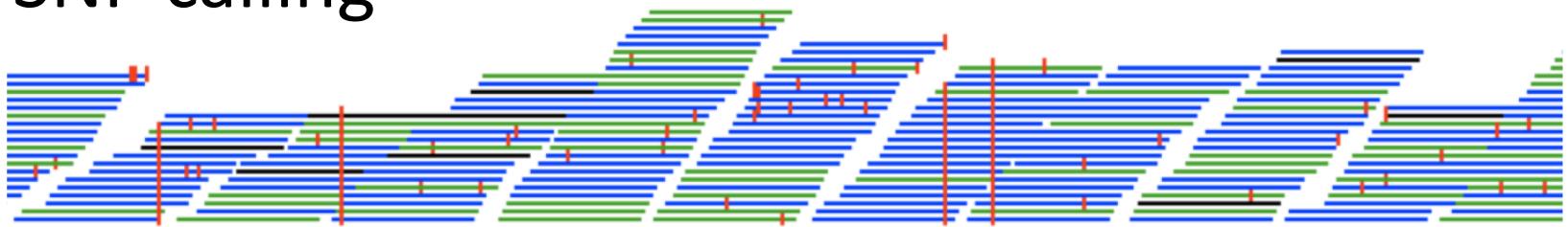


Resequencing and mapping

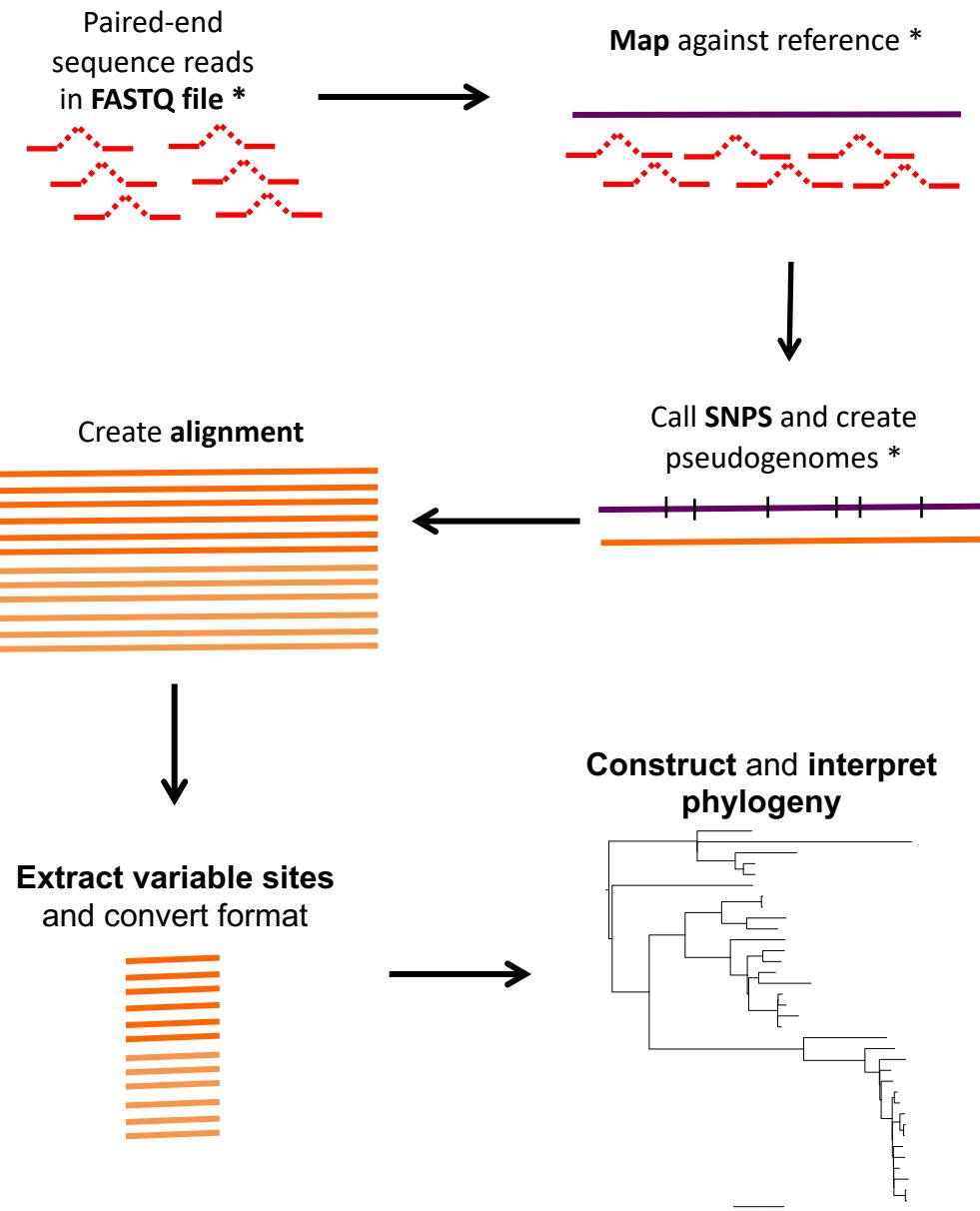
- aims to capture information on:
 - Single Nucleotide Polymorphisms (SNPs),
 - insertions and deletions (indels)
 - Copy Number Variants (CNVs) between variants of the same bacteria.
- As sequences diverge from the reference, mapping becomes progressively less effective

What can you do with mapping?

- Phylogenetics
- SNP calling

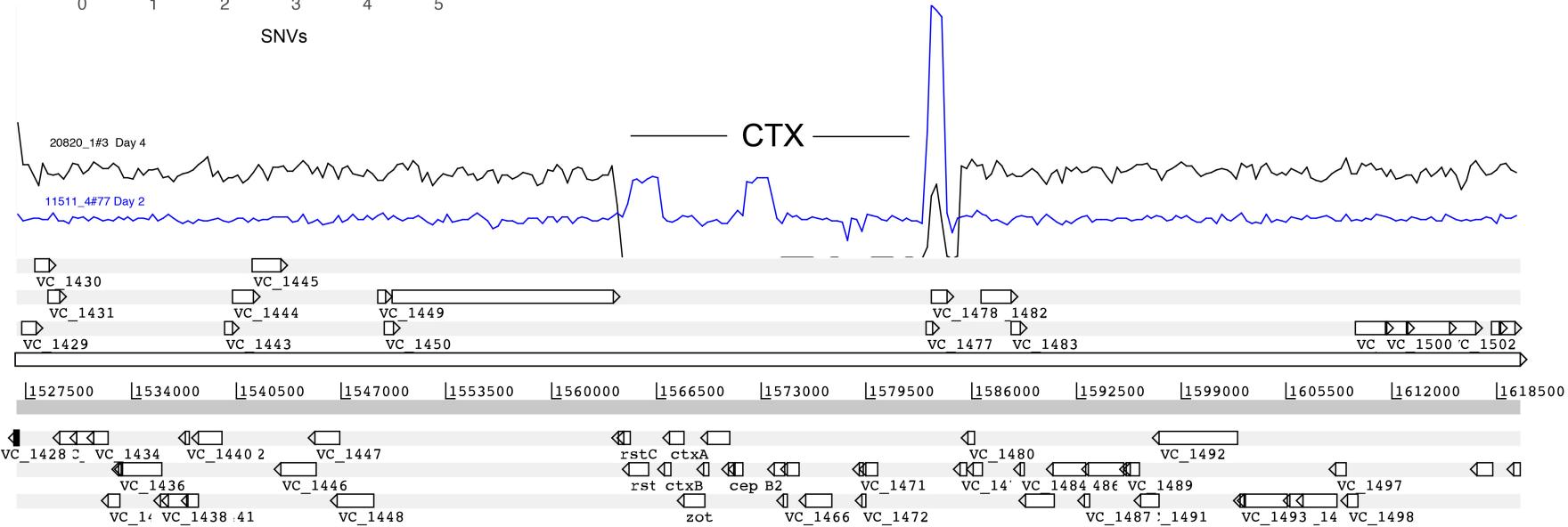
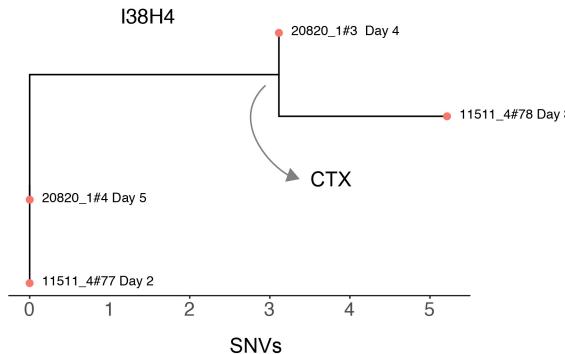


- Looking at copy number
- Looking at presence/absence
- Checking for errors

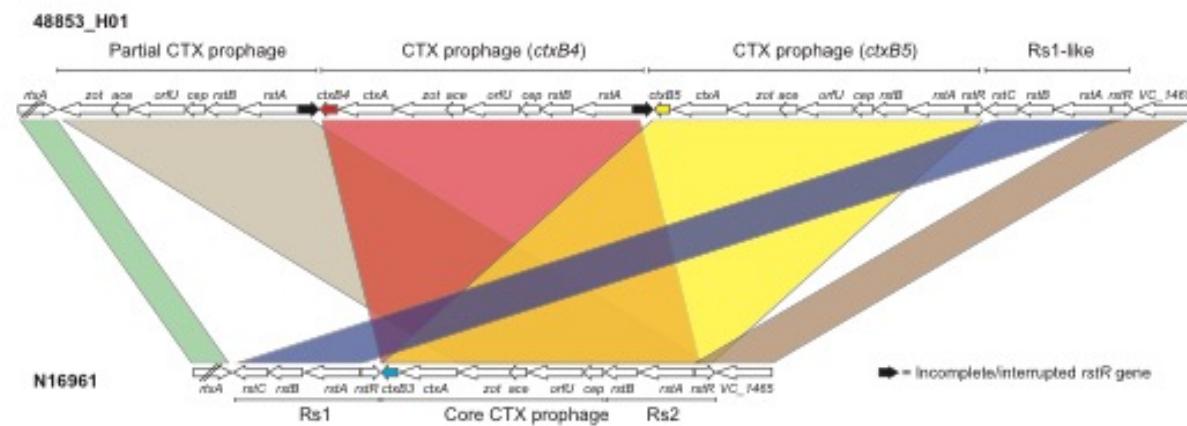
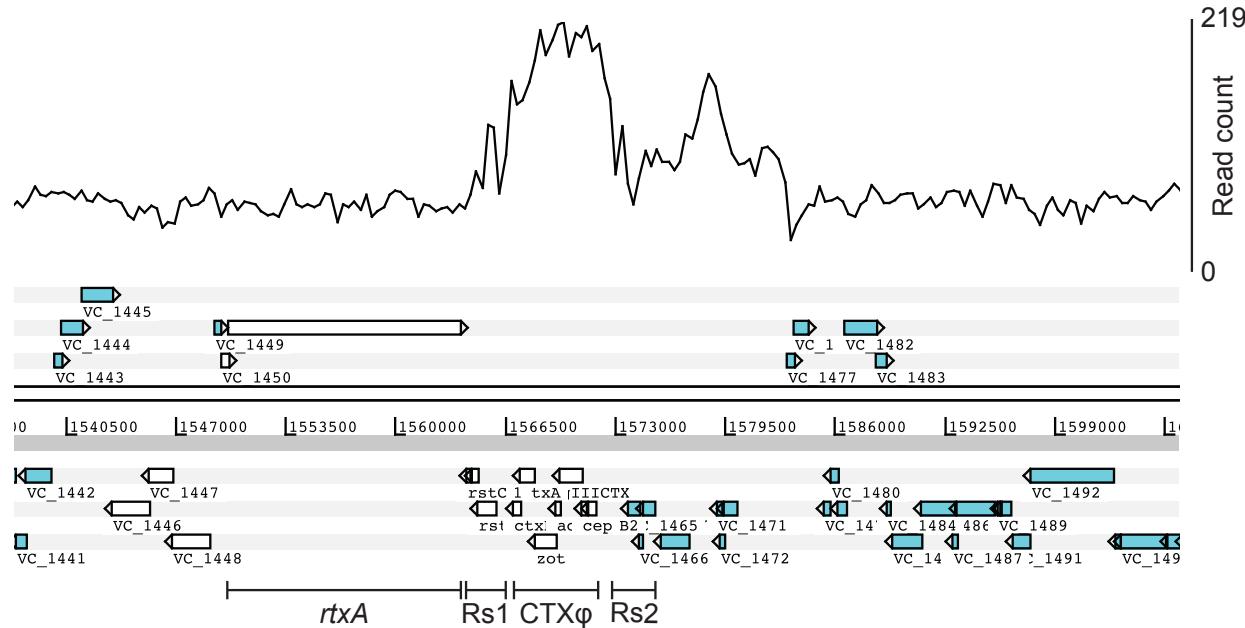


*do for each isolate

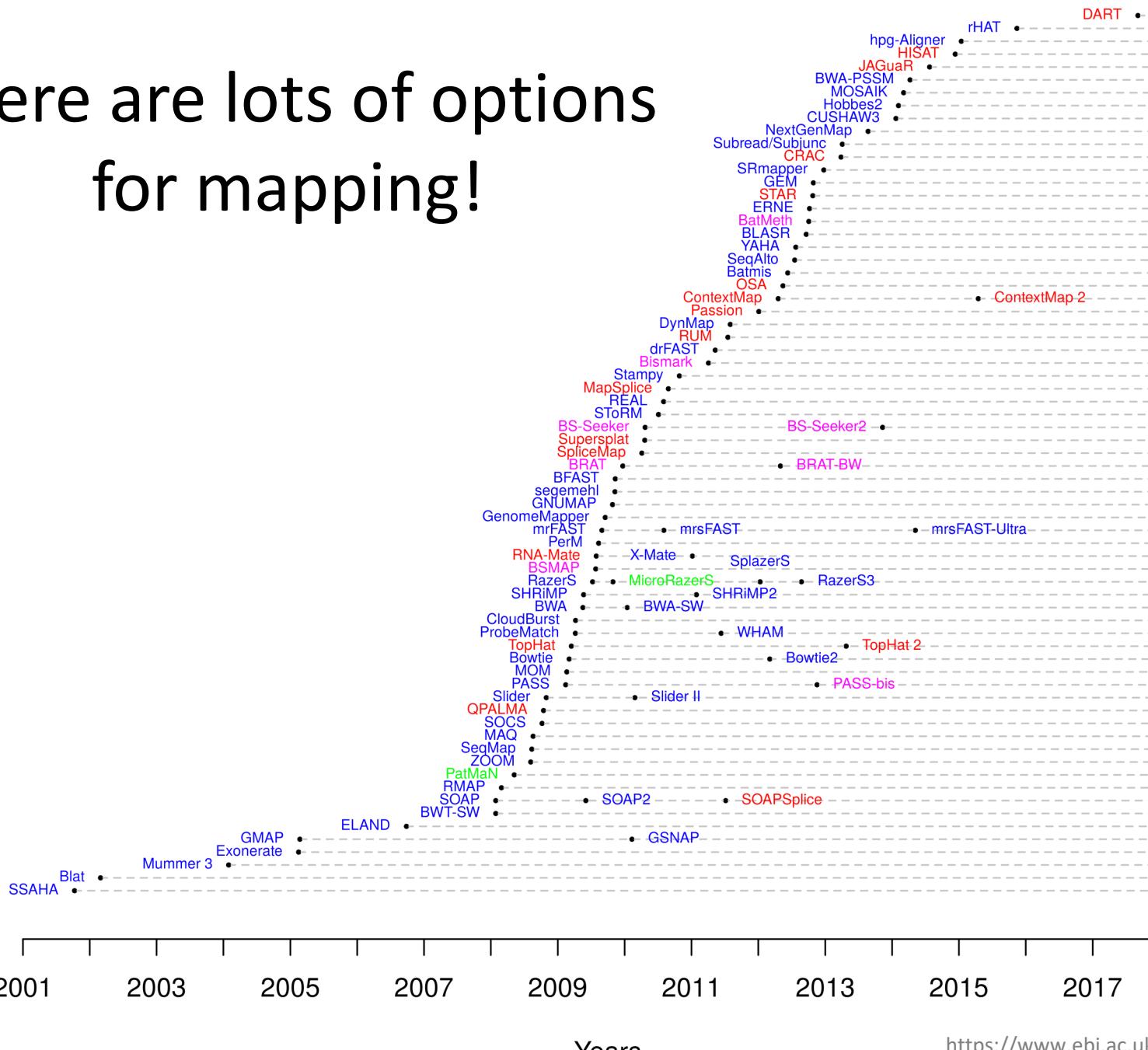
Gene presence / absence



Copy number variation



There are lots of options for mapping!



Comparison of different mappers

Mapper	Data	Availability	Version	O.S.	Number Citations	Seq.Plat.	Input	Output	Min. RL	Max. RL	Mismatches	Indels	Gaps	Align.	Reported	Alignment	Parallel	QA PE	Splicing	Index
BatMeth	Bisulfite	OS	1.03	Linux, Unix	34	I	(C)FAST(A/Q)	Native	35	100	5	N	N	B,U	G	N	Y	N	N Reference	
Batmis	DNA	OS	3.0	Linux, Mac	23	I,So	FASTA/Q	SAM		10	0	N	A,U,S		N	Y	Y	N	N Reference	
BFAST	DNA	OS	0.7.0	Linux, Mac	553	I,So,4, Hel	(C)FAST(A/Q)	SAM TSV	*		Y	Y		B,R,U	G	SM	N	Y	N Reference	
Bismark	Bisulfite	OS	0.7.3	Linux, Mac	887	I	FASTA/Q	SAM	16	10K	Score	Score	N	U		SM	Y	Y	N Reference	
BLASR	DNA	OS	1.4	Linux, Unix	P		FASTA/Q hdf5	SAM TSV	50	100000	0.2	0.2	Y	A,B,R	G L	N	Y	N	De novo Reference	
Blat	DNA	OS	34	Linux, Mac	6252	N	FASTA	TSV BLAST	11	50000	Score	Score	Y	B	L	N	N	N	De novo Reference	
Bowtie	DNA	OS	0.12.7	Linux, Mac, Windows	11207	I,So,4, Sa,P	(C)FAST(A/Q)	SAM TSV	4	1K	Score	Score	N	A,B,R,S	G L	SM	Y	Y	N Reference	
Bowtie2	DNA	OS	2.0beta5	Linux, Mac, Windows	8586	I,4,Ion	FASTA/Q	SAM TSV	4	50000	Score	Score	Y	A,B,R,S	G L	SM	Y	Y	N Reference	
BRAT	Bisulfite	OS	1.2.3	Linux	60	I	FASTA/Q	TSV			Y	0	N			N	N	Y	N Reference	
BRAT-BW	Bisulfite	OS	2.0.1	Linux	53	I	FASTA/Q	TSV	32	*	Y	0	N			N	N	Y	N Reference	
BS-Seeker	Bisulfite	OS		Linux, Mac	193	I	FASTA/Q	TSV			3	0	N	U		SM	Y	N	N	
BS-Seeker2	Bisulfite	OS	2.0.0	Linux, Unix, Mac	107	I	FASTA/Q qseq	SAM BAM	10	200	Score	Score	Y	B,U,S	G L	SM	N	Y	N Reference	
BSMAP	Bisulfite	OS	2.73	Linux, Unix, Mac	347	I	FASTA/Q SAM/BAM	SAM BAM	20	144	15	1	N	B,R,U	G	SM	N	Y	N Reference	
BWA	DNA	OS	0.6.2	Linux, Mac, Windows	13341	I,So,4, Sa,P	FASTQ/Q PSSM	SAM BAM	4	200	Y	8	Y	R,S	G	SM	Y	Y	N Reference	
BWA-PSSM	DNA	OS	0.5.11	Linux	26	I,Hel	FASTQ/Q PSSM	SAM BAM	4	200	30	10	N	B,R,U,S	G	SM	Y	Y	N Reference	
BWA-SW	DNA	OS	0.6.2	Linux, Mac, Windows	3494	I,4,Sa,Hel,Ion,P	FASTA/Q	SAM	4	1000K	0.1	0.1	Y	R,S	L	SM	Y	N	Both	
BWT-SW	DNA	OS	20070916	Linux	133	N	FASTA	TSV		1K	Score	Score	Y	A		N	N	N	N Reference	
CLC Mapper	DNA	Com	4		1,I,So,Sa,Ion,P,Hel		FASTA/Q	SAM BAM							N	Y				
CloudBurst	DNA	OS	1.1	Linux, Mac, Windows	650	N	FASTA	TSV		1K	Y	Y	Y	A,B	G	Cloud	N	N	N Reads	
ContextMap	RNA	OS	2.2	Linux, Unix, Mac	22	I,4,So,Sa,Ion,P,Hel	FASTA/Q	SAM	1	5000	20	10	Y	A,B	G	SM	N	Y	Lib or de novo Reference	
ContextMap 2	RNA	OS	2	Windows, Linux, Unix, Mac	9	I,4,So,Ion,P,Hel	FASTA/Q Illumina	SAM BED	20	5000	0.1	10	Y	B	L	No	N	Y	Lib or de novo Reference	
CRAC	RNA	OS	2.0.0	Linux, Unix, Mac	41	I,4,Ion,P	(C)FAST(A/Q) RAW	SAM BAM	50	*	score	score	Y	A,B,U,S	G	SM	N	Y	De novo Both	
CUSHAW3	DNA	OS	v3.0.3	Linux	33	I,So,4,Ion,P	FASTA/Q	SAM	16	4096	score	score	Y	A,B,R,U,S	G L	SM	Y	Y	N Reference	
DART	DNA	OS	1.2.4	Linux	0	I	FASTA/Q	SAM	20		Y	Y	Y	A,U	G	SM	N	Y	De novo Reads	
drFAST	DNA	OS	1.0.0.0	Linux, Unix	23	So	CFasta(A/Q)	SAM DIVET	25	200	Score	N	N	A,B	G	N	N	Y	N Reference	
DynMap	DNA	OS	0.0.20	Linux	2	N	FASTA	TSV	18	8K	5	0	N	B	L	N	N	Y	N Reads	
ELAND	DNA	Com	1	Linux, Unix, Mac	25	I	FASTA		15	150	2	Score	N	B,S	G	N	Y	Y	N	
ERNE	DNA	OS	1 Windows, Linux, Unix, Mac	14	I	FASTA/Q Illumina	SAM BAM Native	15	600	0.1	5	Y	A,R,U,S	G	SM/DM	N	Y	De novo Reference		
Exonerate	DNA	OS	2.2	Linux, Mac	918	N	FASTA	TSV	20	*	Score	Score	Y	B,S	G L	N	N	N	De novo Reference	
GEM	DNA	Bin	1.x	Linux, Mac	260	I, So	FASTA/Q	SAM Counts		*	Y	Y	Y	A,S	G	SM	Y	Y	Lib and de novo Reference	
GenomeMapper	DNA	OS	0.4.3	Linux, Mac	144	I	FASTA/Q	BED TSV	12	2K	10	10	Y	A,B,R	G	SM	N	N	N Reference	
GIMAP	DNA	OS	2012-04-27	Linux, Unix, Mac, Windows	868	I,4,Sa,Hel,Ion,P	FASTA/Q	SAM GFF Native	8	*	Y	Y	Y	B	G L	SM	N	Y	De novo Reference	
GNUMAP	DNA	OS	3.0.2	Linux, Mac	80	I	FASTA/Q Illumina	SAM TSV	16	1K	Score	Score	Y	B	G	SM/DM	Y	Y	De novo Reference	
GSNAP	DNA	OS	2012-04-27	Linux, Unix, Mac, Windows	1156	I,4,Sa,Hel,Ion,P	FASTA/Q	SAM Native	17	250	Y	Y	Y	A,B,U,S	G L	SM	N	Y	Lib and de novo Reference	
HISAT	RNA	OS	1	Windows, Linux, Unix, Mac	480	I	FASTA/Q	SAM	50	*	0.1	0.1	N	A,B,R,U,S	G	SM	Y	Y	Lib or de novo Reference	
HISAT2	DNA	OS	2	Windows, Linux, Unix, Mac	I		FASTA/Q	SAM	50		score	score	N	B	G	SM	Y	Y	Lib or de novo Reference	
Hobbes2	DNA	OS	2.1	Linux	13	N	FASTA/Q	SAM	22	200	0.08	0.08	N	A,U,S	G	N	N	Y	No Reference	
hpg-Aligner	DNA	OS	v2.1.0	Linux	1 I,4,So,4,Sa,Hel,Ion,P		FASTQ	SAM, BAM	10	2000	0.3	0.3	Yes	A,B	G	N	Y	Y	Lib and de novo Reference	
JAGUAR	RNA	OS	2.1	Linux, Unix	15	I	FASTQ	SAM BAM	50	300	Y	Y	Y	B	G	N	Y	Y	Lib Reference	
MapReads	DNA	OS	2.4.1	Linux, Mac, Windows	0	So	FASTA/Q	TSV	10	120	Score	0	N	S		N	Y	N	N Reference	
MapSplice	DNA	OS	1.15.2	Linux	610	I	FASTA/Q	SAM BED	8	500	3	Y	Y	B	G	SM	N	Y	De novo	
MAQ	DNA	OS	0.7.1	Linux, Mac	2592	I,So	(C)FAST(A/Q)	TSV	8	63	Y	Y	N			N	Y	Y	N Reads	
Masai	DNA	OS	0.4	Windows, Linux, Mac	1	I, Ion	FASTA/Q	SAM TSV	20	32678	32	32	N	A, B, U	G	N	N	Y	N Both	
MicroRazerS	mRNA	OS	0.1	Linux	40	N	FASTA	SAM TSV	10	*	Score	0	N	S		N	Y	N	N Reference	
MIRA	DNA	OS	3	Linux, Unix	48	I,4,Sa,Ion,P	FASTA/Q PHD EXP SAM GFF Counts CAF		25	19000	Score	Score	Y	B,R	L	SM	Y	Y	N Both	
MOM	DNA	Bin	0.6	Linux, Mac, Windows	48	I,4	FASTA	TSV			Y	0	N	A	L	SM	N	Y	N Either	
MOSAIK	DNA	OS	2.1	Linux, Unix, Mac, Windows	174	I,4,Sa,Hel,Ion,P	(C)FAST(A/Q)	BAM	15	1000	Y	Y	Y	A,B	G	SM	Y	Y	N Reference	
mrFAST	DNA	OS	2.5.0.1	Linux, Unix	602	I	FASTA/Q	SAM DIVET	25	1000	Score	4	N	A,B	G	N	Y	Y	N Reference	
mrsFAST	DNA	OS	2.4.0.4	Linux, Unix	229	I	FASTA/Q	SAM DIVET	25	100	Score	N	N	A	G	N	Y	Y	N Reference	
mrsFAST-Ultra	DNA	OS	3.3.1	Linux, Mac	28	I	FASTA/Q	SAM DIVET	8	500	Score	N	N	A,B,S	G	SM	Y	Y	N Reference	
Mummer 3	DNA	OS	3.23	Linux, Mac	2446	N	FASTA	TSV	10	*	Y	Y	Y	A,B	G	N	N	N	N Reference	
NextGenMap	DNA	OS	0.4.6	Linux	82	I,4,Ion	(C)FAST(A/Q),SAM,BAM	SAM BAM	13	1000	Score	Score	N	R,S	G L	SM	N	Y	N Reference	
Novoalign(CS)	DNA	Bin	V2.08.03	Linux	0	I,So,4, Hel,Ion	(C)FAST(A/Q)	SAM Native	1	250	Y	Y	Y	A,B,R,U	G	SM/DM	Y	Y	Lib de novo Reference	
OSA	RNA	Bin	1.0 <	Windows, Linux, Unix, Mac	54	I,4,Ion	FASTA/Q	SAM BAM	15	8000	*	*	Y	A,B,U	G	SM	Y	Y	De novo Reference	
PASS	DNA	Bin	1.62	Linux, Mac, Windows	142	I,So,4	(C)FAST(A/Q)	SAM GFF3 BLAST	23	1K	Y	Y	Y	A,B	G	SM	Y	Y	N Reference	
PASS-bis	Bisulfite	OS	2.01	Linux	14	I,So,4,Sa	FASTA/Q	SAM GFF Counts	14	2000	Score	N	N	A, B, U	G	SM	Y	Y	N Reference	
Passion	DNA	OS	1.2.0	Linux, Unix	28	I,4,Sa,P	FASTA/Q	BED	1	*	Y	Y	Y	A	G	N	N	N	N Reference	
PatMan	mRNA	OS	1.2.2	Linux, Mac	140	N	FASTA	TSV	20	128	9	0	Y	A,U	G	DM	Y	Y	N Reads	
PerfM	DNA	OS	0.4.0	Linux, Unix, Mac, Windows	113	I,So	(C)FAST(A/Q)	SAM TSV	36	50	3	Y	N	A,B	G	N	N	N	N Reference	
ProbeMatch	DNA	OS		Linux, Mac	4	I,4,Sa	FASTA	ELAND			Y	Y	Y	B	L	N	Y	Y	Lib and de novo	
QPALMA	DNA	OS	0.9.2	Linux, Mac	169	I,4	Specific	TSV			Y	Y	Y	B	G	N	Y	Y	N Reference	
RazerS	DNA	OS	1.2	Linux, Mac, Windows	165	I,4	FASTQ	TSV ELAND	11	*	Score	Score	Y	A,B,U,S	G	SM	N	Y	N Reads	
RazerS3	DNA	OS	3.1	Windows, Linux, Mac	81	I	FASTA/Q	SAM TSV GFF	11	*	0.5	Y	N	A,B,U,S	G	SM	N	Y	N Reference	
REAL	DNA	OS	0.0.28	Linux	32	I	FASTA/Q	TSV	4	*	Score	N	N	B,U	G	SM	Y	Y	N Reference	

Good general aligners

bwa

→

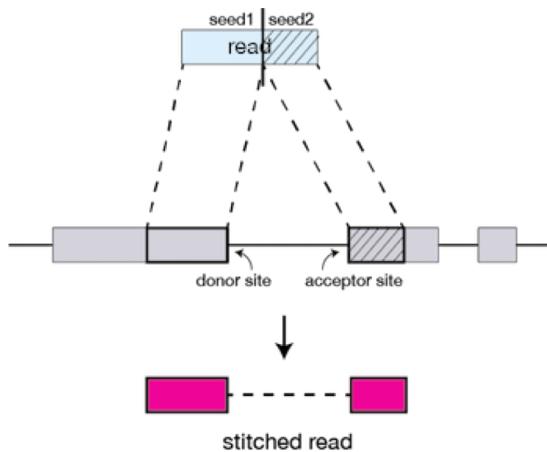
bowtie2

Fast, sensitive and
easy to use!

Splice-aware aligners for RNA-seq

STAR

HISAT2



Mapping sequencing reads-Workflow

Fastq files

Mapping script

mapping (bwa-mem)

SAM file

sam/bam conversion
(samtools)

BAM file

snp calling
(samtools)

BCF file

1

Artemis

2

3

running
mapping and
snp calling

SNPs, INDELs



You must think carefully about what reference genome you want to use!

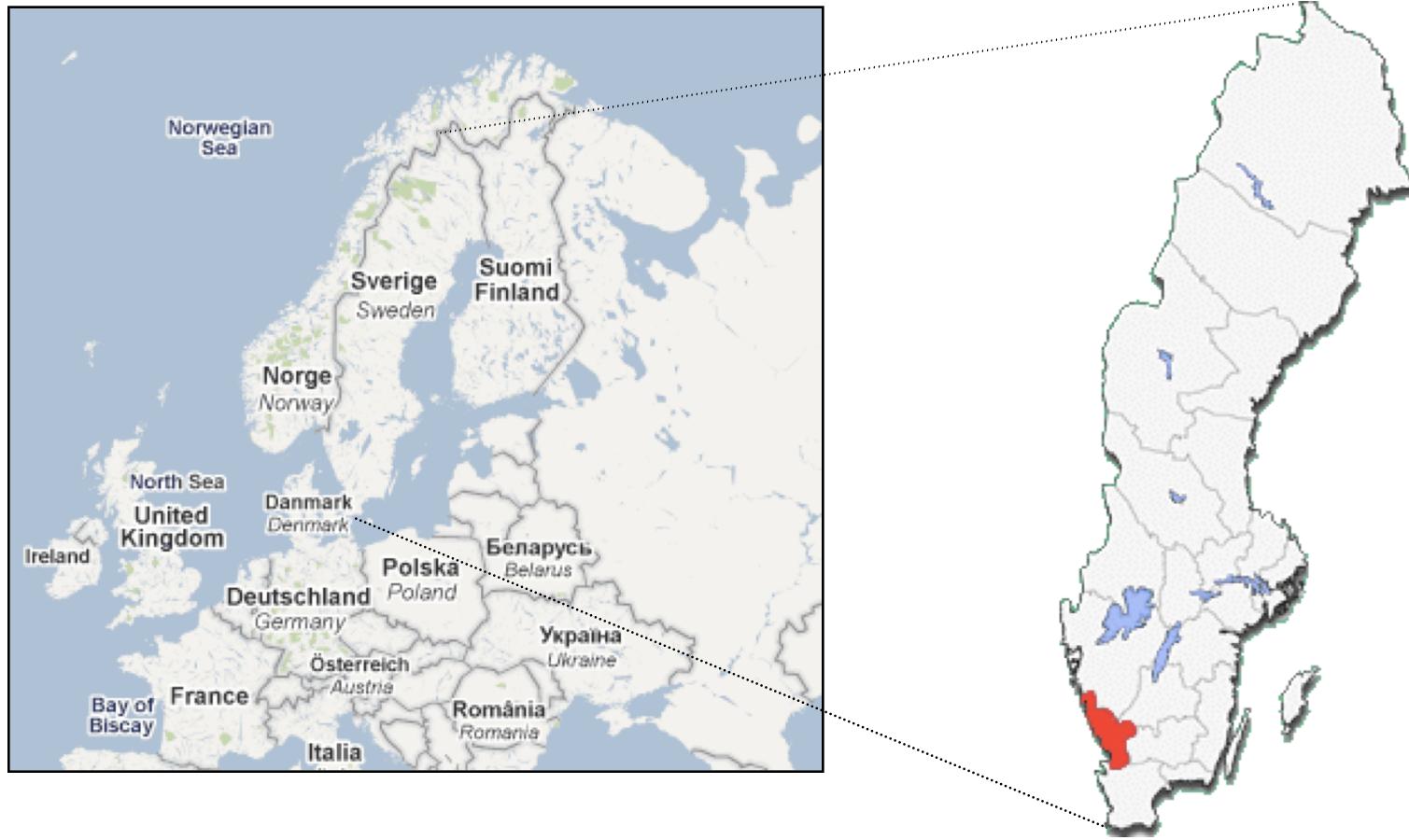
You won't find things in your sample that are not in the reference!

The Swedish Story



- Prior to 2006, C. *trachomatis* in Sweden was following the same pattern as in the UK
- In 2006, across Sweden there was a reported 2% drop in cases
- Initially thought this was a good thing

The County of Halland, Sweden



In Halland County there was a 25% decrease in diagnosed *C. trachomatis* infections from Nov 05 - Aug 06

The Swedish story

- It was noticed that counties using the NAATs (Abbott / Roche) diagnostic system showed a drop in *C. trachomatis* cases in 2006
- Halland County showed the greatest drop
- Counties using other diagnostic methods (Beckton Dickinson) still showed an increase in cases, in line with that of previous years
- The obvious conclusion was that the NAATs was missing a subset of infections
- Halland County Like most places had dropped the dual site PCR for the cheaper and apparently just as accurate single site plasmid PCR.

Sequencing pSW2 (nvCT)

