

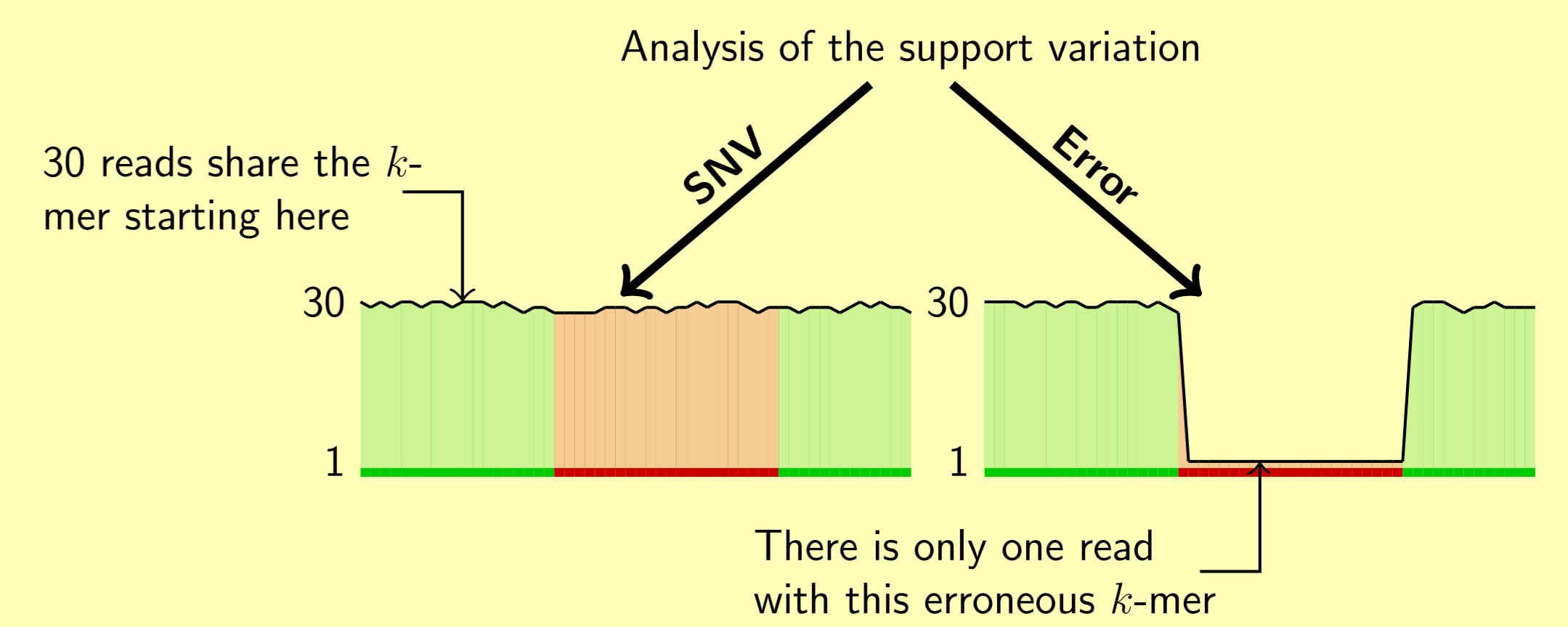
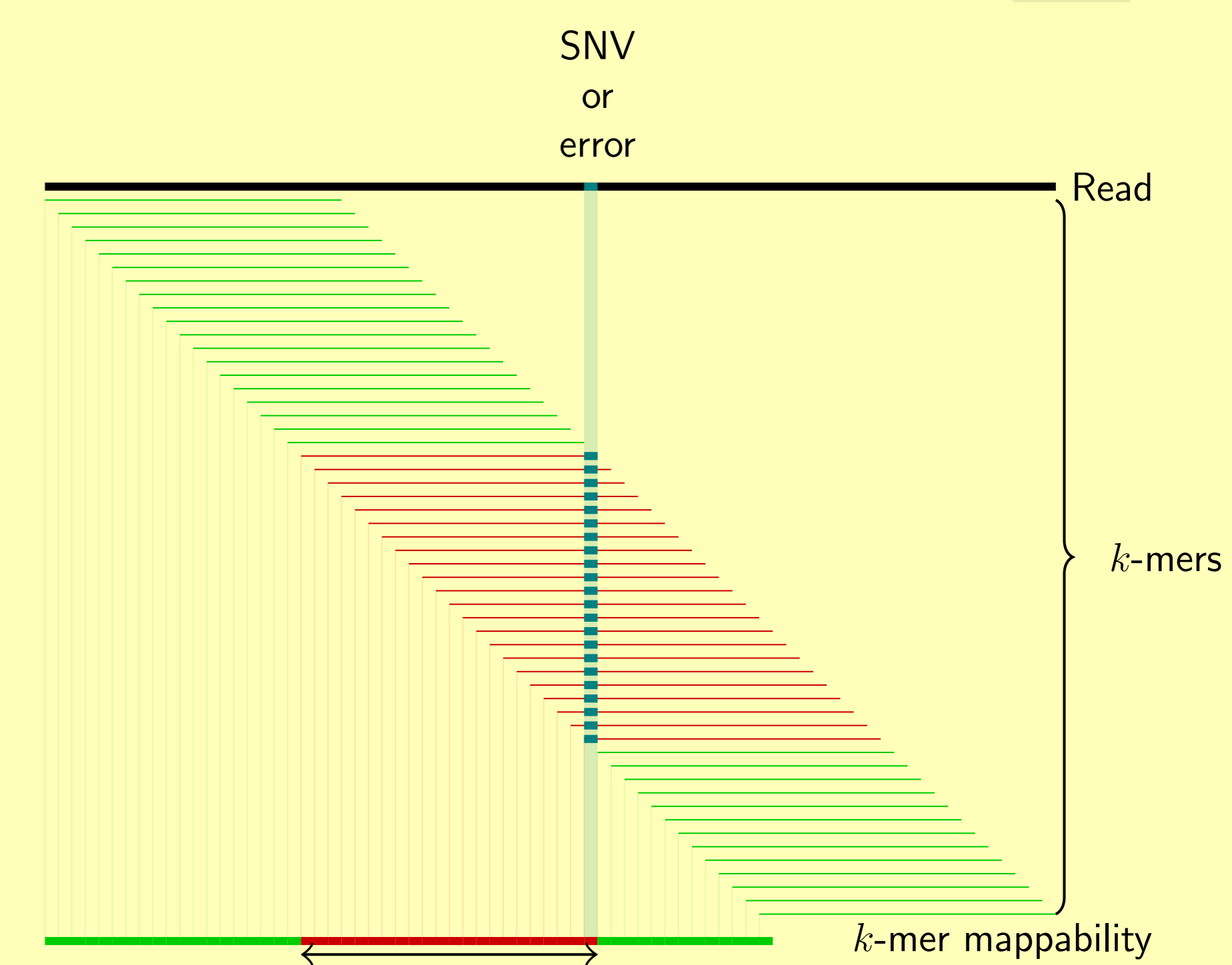
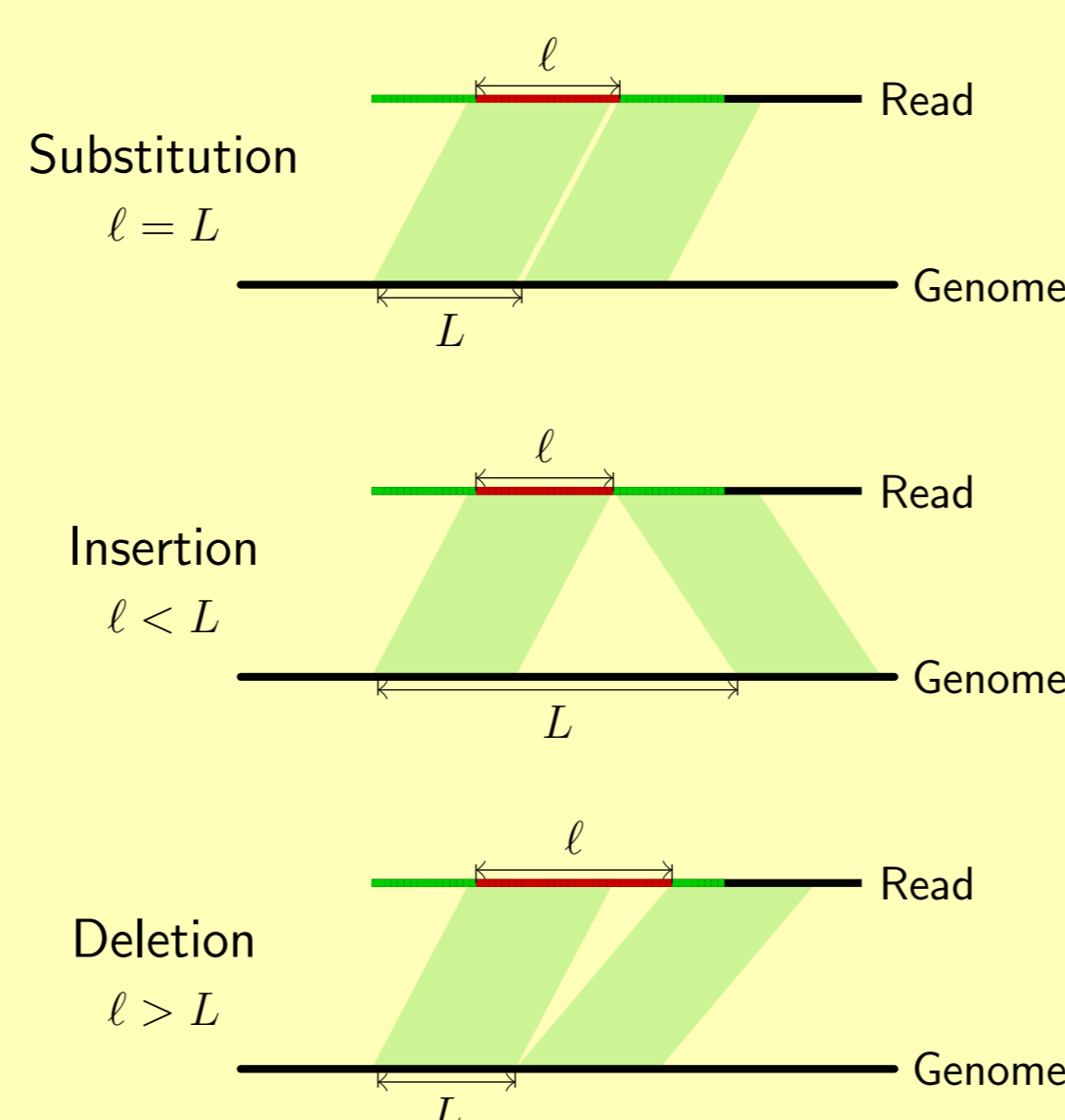
Abstract

NGS techniques are intensively used to investigate transcriptome complexity using RNA-sequencing assays. Huge read sets are then compared to a reference genome to determine transcribed exons, and unmapped reads are further analyzed for predicting splice junctions. Subsequent analyses seek to reconstruct standard and alternative transcripts. However, transcript inference can only succeed if reads are mapped correctly and splice junctions predicted accurately. By comparing current solutions we gathered evidence suggesting that read processing was far from fully completing this task. Here, we propose a novel way of analyzing reads that integrates genomic locations and local coverage to distinguish sequence errors from biological mutations, and to infer directly splice junctions within a single read. An evaluation of our program CRAC shows that it improves mapping sensitivity up to 30% compared to state of the art solutions, while being highly specific. Such a sensitivity gain can impact the user's ability to detect rare mutations or splicing variants. Moreover, results indicate that CRAC predicts with bp precision the splice junctions of reads sequenced from non colinear, chimeric RNAs with [40, 60]% sensitivity and > 90% specificity, a unique feature to our knowledge. Finally, this integrated read analysis strategy may broaden the scope of the transcriptome that is amenable to discovery using RNA-sequencing approaches.

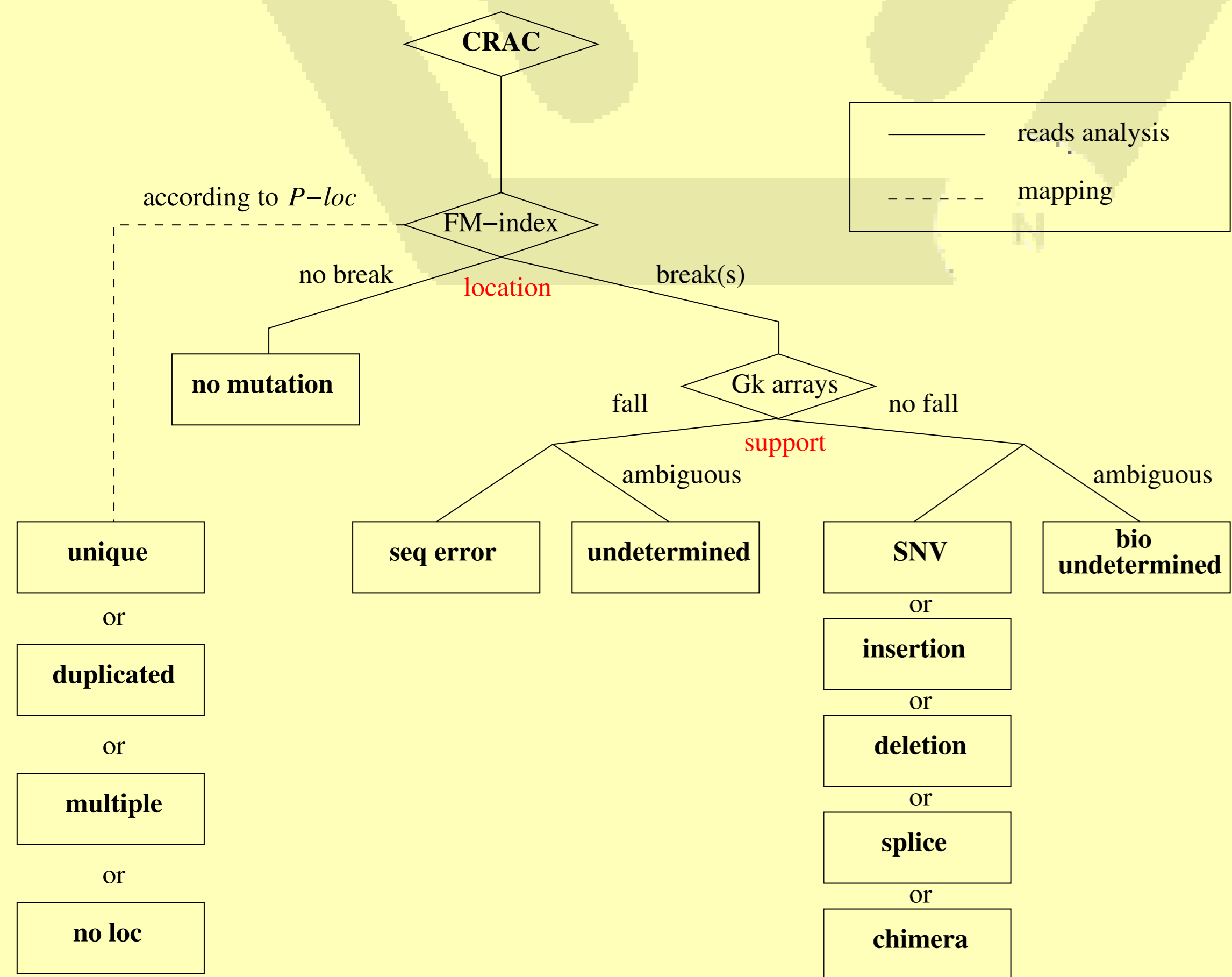
Algorithm

CRAC proceeds each read in turn. For each, it monitors two "signals" that vary with the position in the read sequence of length m . For this, it considers the k -mer starting at every position (ie. $m - k + 1$ possible k -mers) in the read and registers:

1. the exact mappability of the k -mer on the reference genome, its matching locations and their number,
2. the k -mer support, which we define as the number of reads sharing this k -mer (ie, the exact same k -mer sequence matches a k -mer from another read). The support value has a minimum value of one since the k -mer exists at least in the current read.



CRAC strategy is to analyse jointly k -mer support and k -mer mappability to predict in a single analysis sequencing errors, as well as potential genetic variations, splice junctions, or chimeras



Results

Biological validations on private AML library:

- ~ 40 millions of unoriented 100 bp reads
- Detection of **511** chimeras

Comparative evaluation of mapping sensitivity and precision

Tool	75bp		200bp	
	Sensitivity	Precision	Sensitivity	Precision
Bowtie	75.42	99.59	55.72	99.81
BWA	79.29	99.13	68.66	96.86
CRAC	94.51	99.72	96.02	99.92
GASSST	70.73	99.09	59.43	97.86
GSNAP	94.62	99.88	84.84	99.28
SOAP2	77.6	99.52	56.08	99.78

Comparative evaluation splice junction prediction tools

Tool	75bp		200bp	
	Sensitivity	Precision	Sensitivity	Precision
CRAC	79.43	99.5	86.02	99.18
GSNAP	84.17	97.03	72.94	97.09
MapSplice	79.89	97.68	84.72	98.82
TopHat	84.96	89.59	54.07	94.69

Comparative evaluation of chimeric RNA prediction tools

Tool	75bp		200bp	
	Sensitivity	Precision	Sensitivity	Precision
CRAC	53.89	93.84	64.86	90.18
MapSplice	2.33	0	2.63	0.01
TopHatFusion	32.73	42.02		
TopHatFusionPost	12.26	97.22		

ECHANTILLES	CT	TM	integrate	CT	TM	integrate	CT	TM	integrate	CT	TM	integrate
OS110089	Applied	27,250,96	92,980,03	Applied	21,205,47	79,366,02	Applied	24,008,00	82,281,08	Applied	40	88,824,03
OM100069	Promega	28,002,02	93,082,06	Promega	28,002,02	93,082,06	Promega	28,002,02	93,082,06	Promega	40	88,824,03
OM110171	Applied	40,500,00	92,840,00	Applied	40,500,00	92,840,00	Applied	40,500,00	92,840,00	Applied	40	88,824,03
OM110448	Promega	27,250,96	92,980,03	Promega	27,250,96	92,980,03	Promega	27,250,96	92,980,03	Promega	40	88,824,03
R2020	Applied	38,52	82,75	Applied	38,52	82,75	Applied	38,52	82,75	Applied	40	88,824,03
OS100380	Applied	40,00	92,840,00	Applied	40,00	92,840,00	Applied	40,00	92,840,00	Applied	40	88,824,03
OM110564	Promega	27,250,96	92,980,03	Promega	27,250,96	92,980,03	Promega	27,250,96	92,980,03	Promega	40	88,824,03
OM110522	Applied	37,02	82,07	Applied	37,02	82,07	Applied	37,02	82,07	Applied	40	88,824,03
OM110520	Promega	37,02	82,07	Promega	37,02	82,07	Promega	37,02	82,07	Promega	40	88,824,03
OM110424	Applied	37,02	82,07	Applied	37,02	82,07	Applied	37,02	82,07	Applied	40	88,824,03
OS110518	Promega	37,02	82,07	Promega	37,02	82,07	Promega	37,02	82,07	Promega	40	88,824,03
OM100590	Applied	37,02	82,07	Applied	37,02	82,07	Applied	37,02	82,07	Applied	40	88,824,03
OM100520	Promega	37,02	82,07	Promega	37,02	82,07	Promega	37,02	82,07	Promega	40	88,824,03
OM100511	Applied	37,02	82,07	Applied	37,02	82,07	Applied	37,02	82,07	Applied	40	88,824,03
OM100503	Promega	37,02	82,07	Promega	37,02	82,07	Promega	37,02	82,07	Promega	40	88,824,03
Cell	Applied	37,02	82,07	Applied	37,02	82,07	Applied	37,02	82,07	Applied	40	88,824,03
BLANC	Promega	37,02	82,07	Promega	37,02	82,07	Promega	37,02	82,07	Promega	40	88,824,03

Conclusions

Highlights:

- Low false positive rate
- Between 60 and 70 % of causes are found (mutations not found are due to a low coverage)
- junctions: more sensitive and more specific than GSNAP, MapSplice, and TopHat

Futur works:

- Transcripts reconstruction (assembly)
- Clinical markers for prognostic and diagnostic
- Chimera variants in myeloid leukemia (normal karyotype)

CRAC is particularly suitable for the data of the future: more massive and longer