

Challenges and Pitfalls in the Evaluation of Structural Variants Callers

Luca Denti¹, Thomas Krannich², Tomáš Vinař¹, Rayan Chikhi³, Paola Bonizzoni⁴, Broňa Brejová¹, Fereydoon Hormozdiari⁵

¹ Comenius University in Bratislava, Slovakia; ² German Cancer Research Center, Germany; ³ Institut Pasteur, France

⁴ University of Milano-Bicocca, Italy; ⁵ University of California-Davis, USA



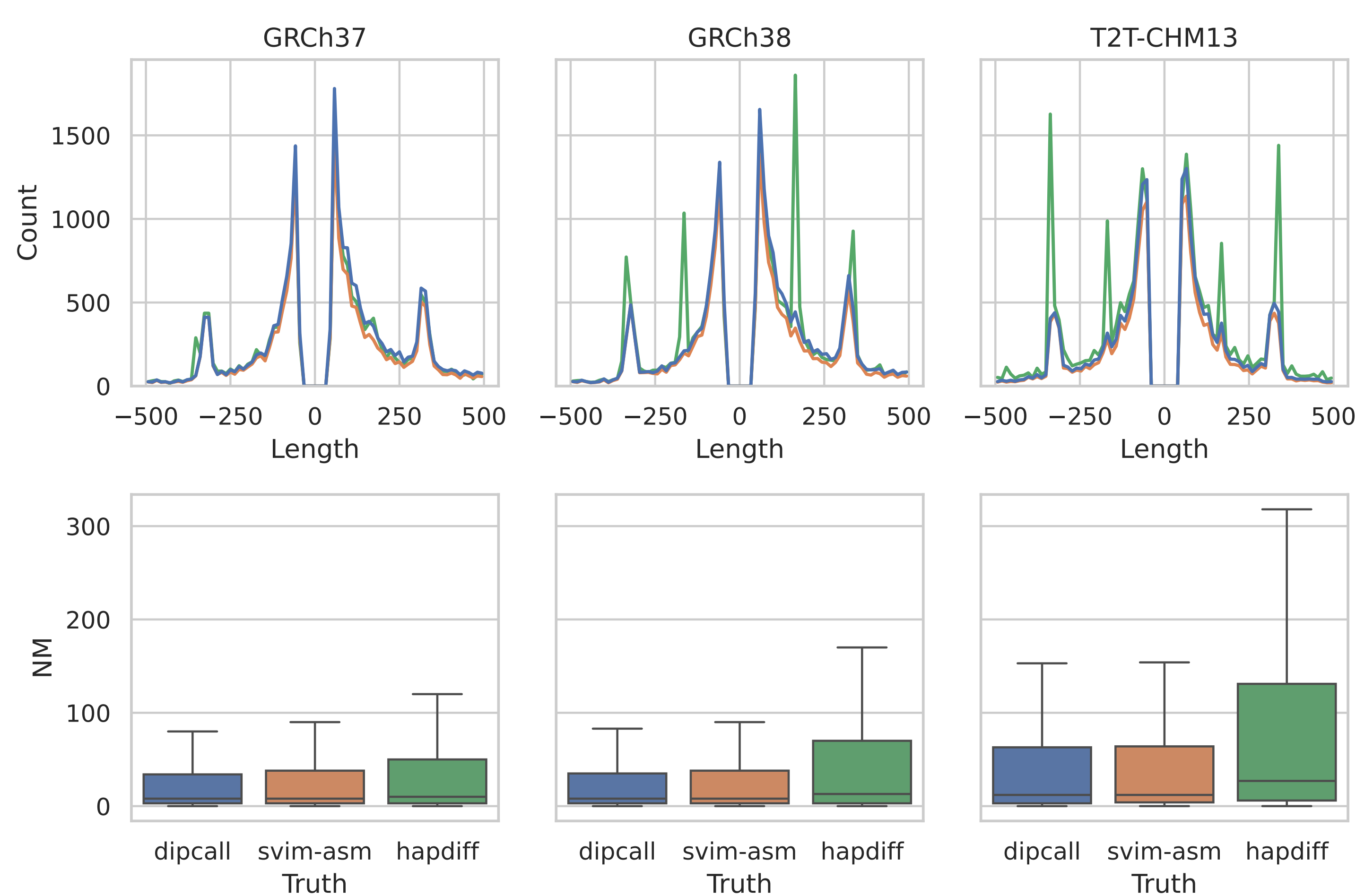
Motivation

Structural variants (SVs) are medium- and large-scale genomic alterations that shape phenotypic diversity and disease risk. Numerous methods have been proposed for the discovery of SVs, but their benchmarking has been inconsistent across studies, often resulting in *contradictory findings*. Primary sources of conflicting results lie in the inconsistent selection of reference genome version and ground truths, e.g., callsets curated by consortia and ad hoc callsets created from high-quality assemblies. Best practices and optimal evaluation strategies remain undefined and benchmarking results often vary, leaving the current performance assessment landscape *difficult to interpret*. More importantly, this has created *ambiguity* about the true capabilities of leading methods to accurately detect SVs, with far-reaching implications for routine genome analysis.

SV calling from high-quality *de novo* assemblies

HG002 (GIAB Q100), HG002 (HPRC), NA12878 (Platinum Pedigree)

(i) Leading methods produced different callsets, even when applied to the same assembly.



(ii) Similarity decreases as the completeness of the reference genome increases.

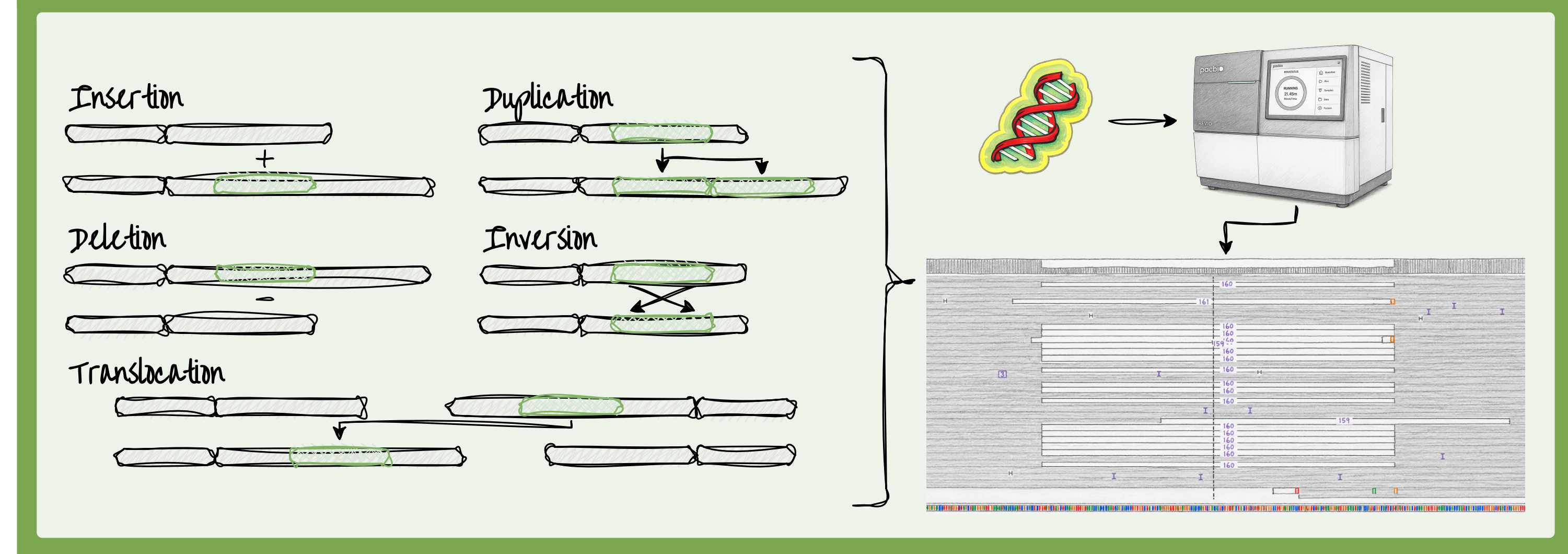
	GRCh37	GRCh38	T2T-CHM13
dipcall	1.0	0.91	0.83
SVIM-asm	0.91	1.0	0.83
hapdiff	0.83	0.83	1.0

(iii) Differing computational strategies applied to address the biological complexity of certain genomic loci are one of the primary sources of these discrepancies:

- alignment choices
- signal clustering (within and across haplotypes)



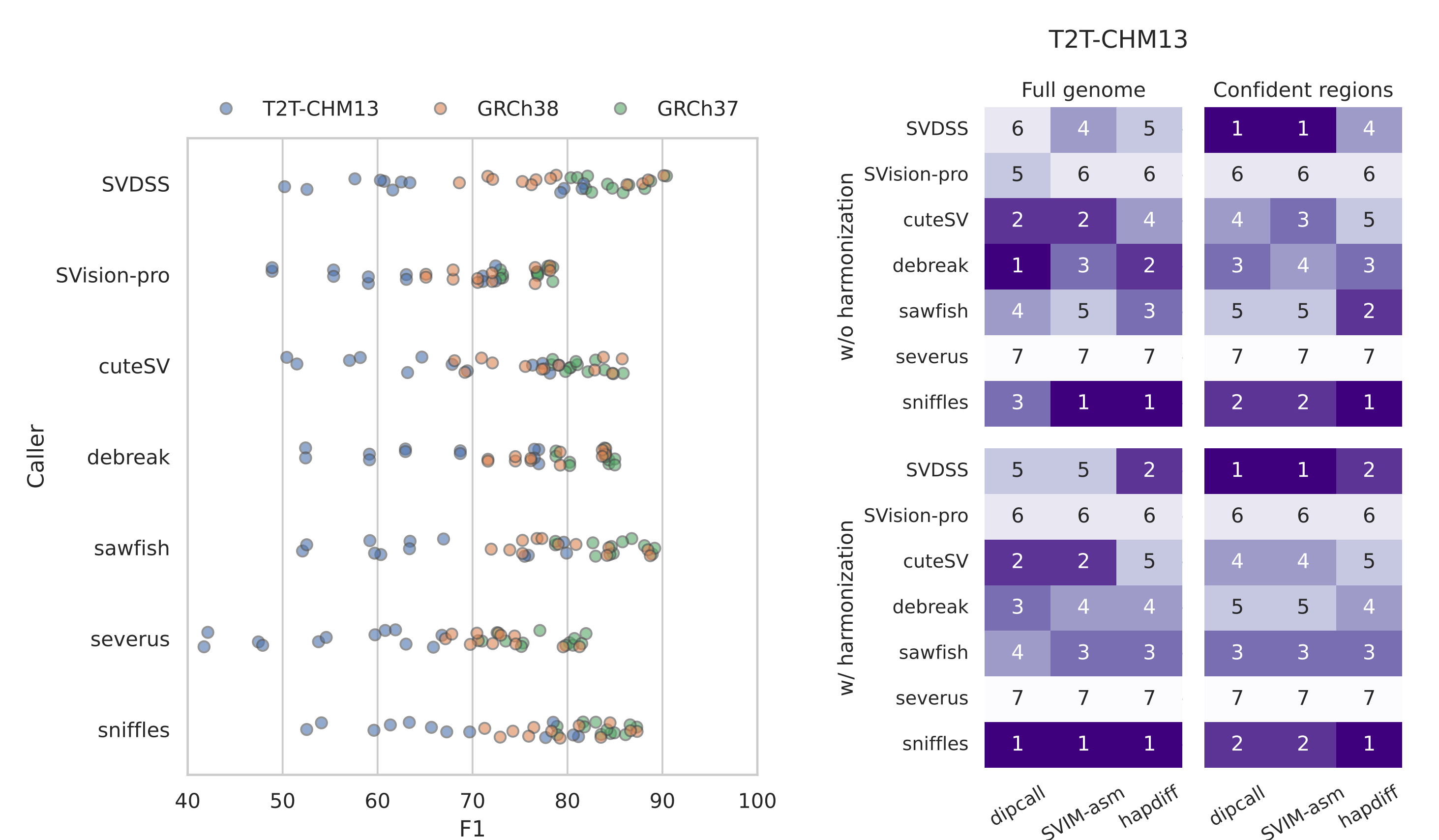
Background



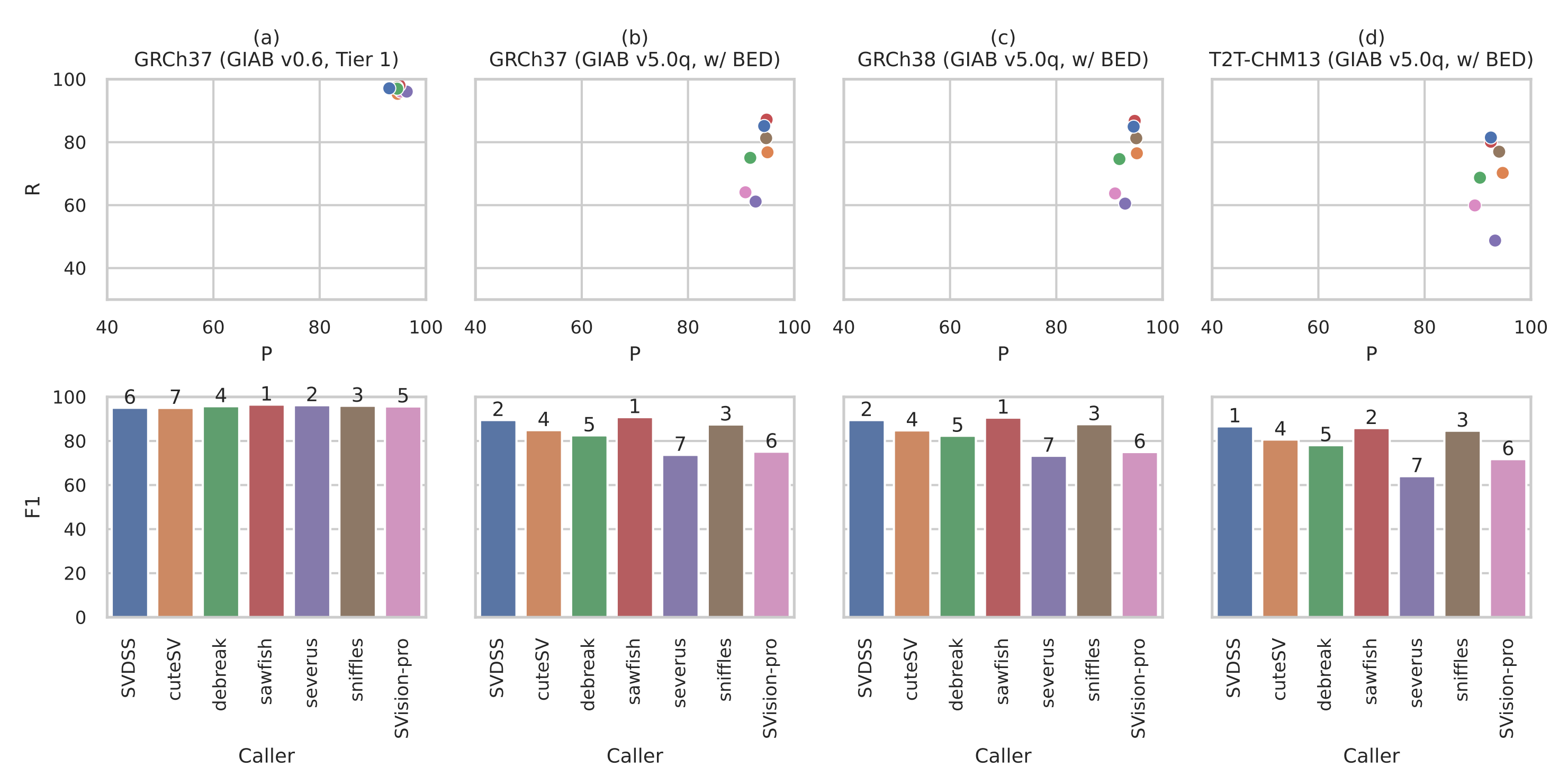
SV calling from PacBio HiFi reads

HG002 (36x, HPRC), NA12878 (100x, Platinum Pedigree) - Sequel II

(i) When evaluated against assembly-based callsets, F1-scores and rankings of callers vary depending on the choice of ground truth, reference genome, and benchmark parameters.



(ii) Similar results can be observed when considering curated SV callsets. Remarkably, the inclusion of more (challenging) SVs located in complex genomic regions substantially affects SV callers performances.



Discussion

- Common user-defined choices cause overlooked and unexpected divergence in performance evaluations.
- Performance results are not readily transferable between ground truths without careful consideration of their subtle differences.
- All long-read SV callers struggle more in complex genomic regions. Both calling and benchmarking are more challenging (due to inconsistent SV representation).
- Although some callers perform well across our benchmarks, no tool was best in every tested scenario.
- Due to the complexity of SV comparison, evaluating methods against a single “ground truth” is insufficient.