



Scalable subclonal reconstruction of cancer cells in DNA sequencing data using a penalized likelihood model

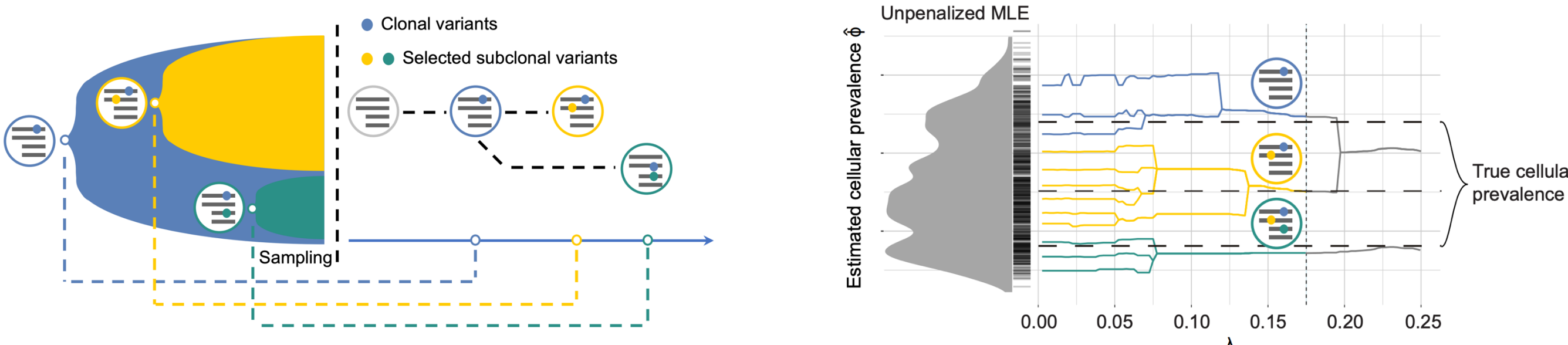
Yujie Jiang^{1,2,3}, Matthew D Montierth^{1,3,4}, Yu Ding¹, Kaixian Yu⁴, Quang Tran¹, Aaron Wu¹, Ruonan Li¹, Shuangxi Ji¹, Xiaoqian Liu^{1,5}, Seung Jun Shin⁶, Shaolong Cao¹, Yuxin Tang⁷, Tom Lesluyes⁸, Marek Kimmel², Jennifer R. Wang⁹, Maxime Tarabichi¹⁰, Hongtu Zhu¹¹, Peter Van Looy^{8,12,13}, Wenyi Wang^{1*}

¹Department of Bioinformatics and Computational Biology, The University of Texas M.D. Anderson Cancer Center, Houston, TX 77030 ²Department of Statistics, Rice University, Houston, TX 77005 ³Graduate program in Quantitative Computational Biology, Baylor College of Medicine, Houston, TX 77030 ⁴Department of Biostatistics, The University of Texas M.D. Anderson Cancer Center, Houston, TX 77030 ⁵Current address: Department of Statistics, The University of California, Riverside, CA 92507 ⁶Department of Statistics, Korea University, Seoul, South Korea, 02841 ⁷Department of Computer Science, Rice University, Houston, TX 77005 ⁸Cancer Genomics Group, The Francis Crick Institute, London, UK ⁹Department of Head and Neck Surgery, The University of Texas M.D. Anderson Cancer Center, Houston, TX 77030 ¹⁰RIBHM-J.E. Dumont, Brussels, Belgium, B1070 ¹¹Department of Biostatistics, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599 ¹²Department of Genetics, The University of Texas M.D. Anderson Cancer Center, Houston, TX 77030 ¹³Department of Genomic Medicine, The University of Texas M.D. Anderson Cancer Center, Houston, TX 77030 *authors contributed equally. *Correspondence: W.Wang7@mdanderson.org



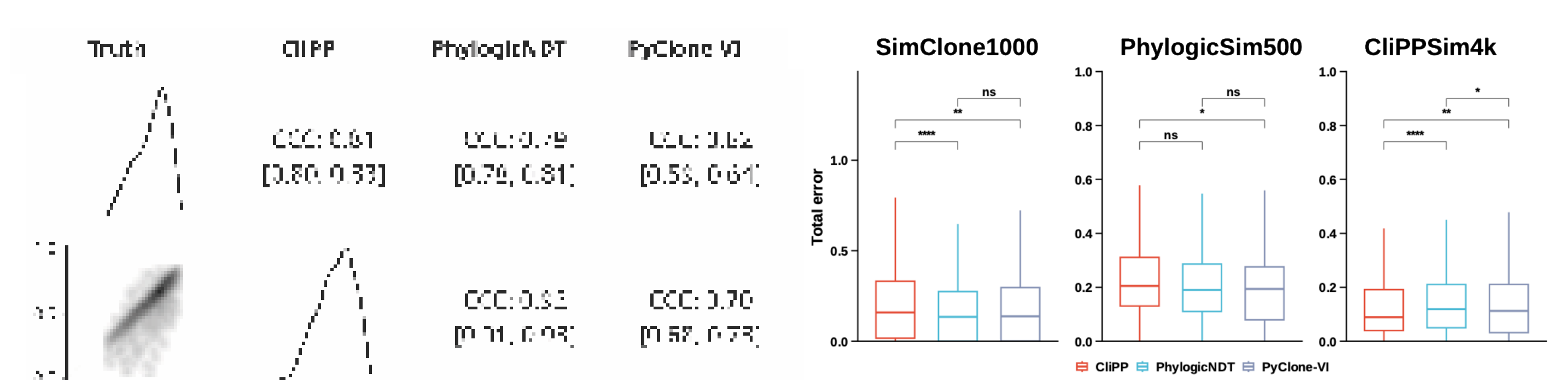
CliPP: Clonal structure identification through Pairwise Penalization

CliPP¹ estimates mutation-specific cellular prevalence from SNV read counts after adjusting for tumor purity and allele-specific copy number. Pairwise SCAD² penalization fuses mutations with similar cellular prevalence, producing clonal and subclonal mutation clusters without pre-specifying the number of clones.



Input: SNV read counts + tumor purity + allele-specific copy number
Output: mutation clusters + cellular-prevalence estimates + clonal/subclonal labels

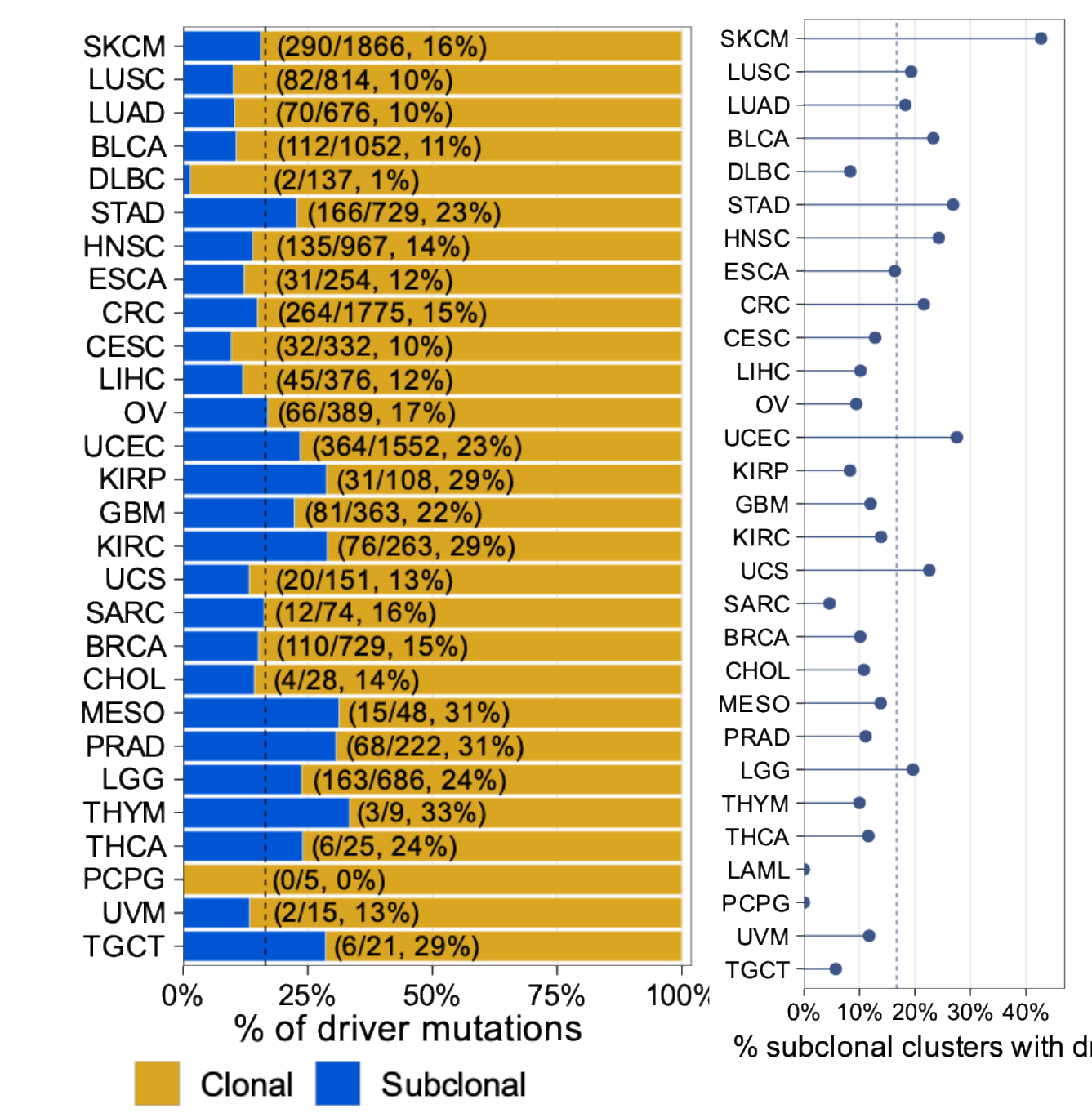
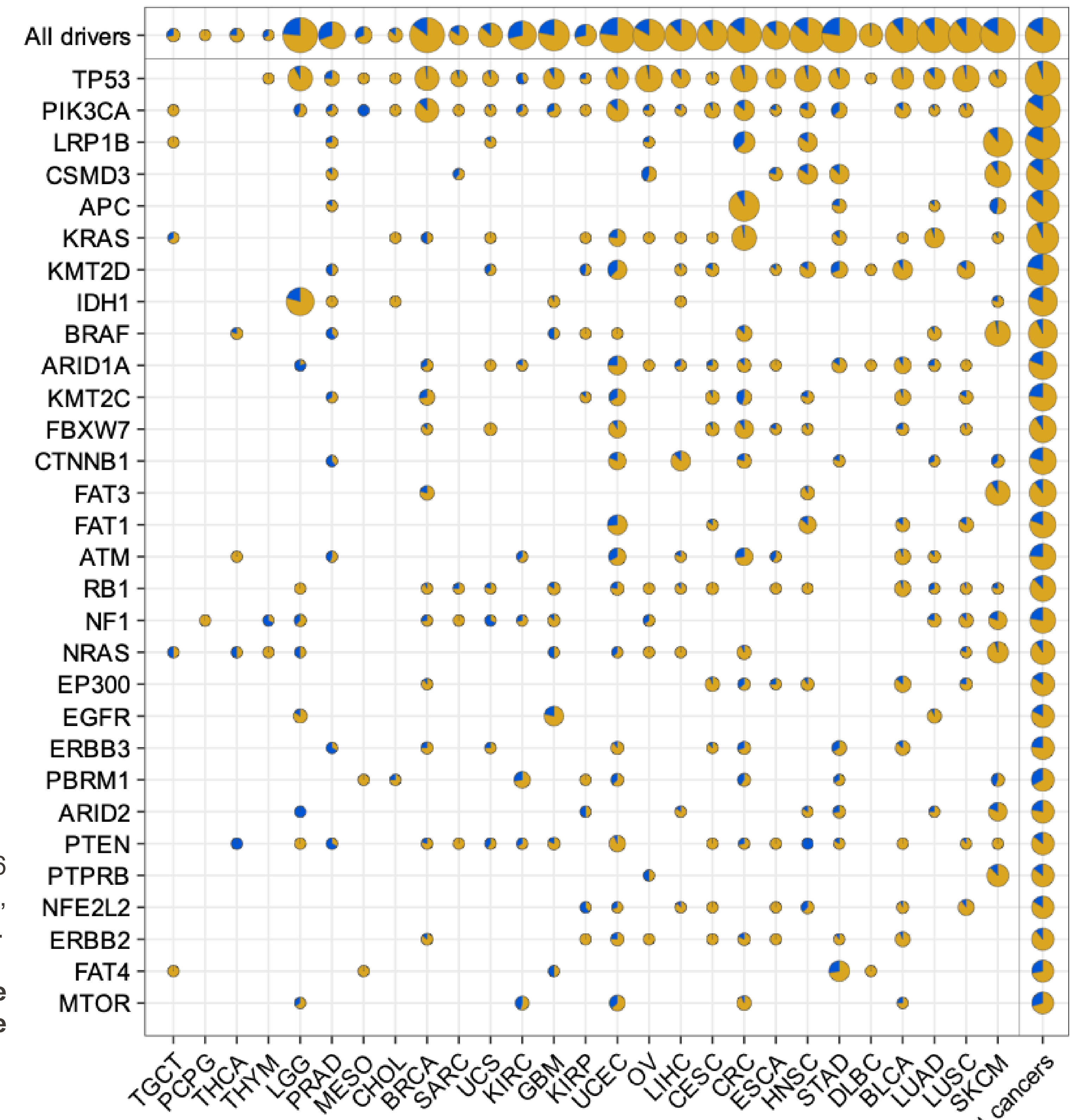
CliPP is fast and accurate in simulated and real data



CliPP was benchmarked on 5,306 simulated tumors from CliPPSim4k, PhylogSim500^{3,4}, and SimClone1000^{3,4}.

CliPP provides accuracy comparable to leading methods with runtime suitable for large-cohort workflows.

Pan-cancer driver clonality landscape



CliPP annotated 13,666 high-confidence driver mutations in 274 genes across 4,218 tumors; 28 cancer types are shown in the final driver-clonality panels.

2,256 / 13,666 driver mutations were subclonal (16.5%).

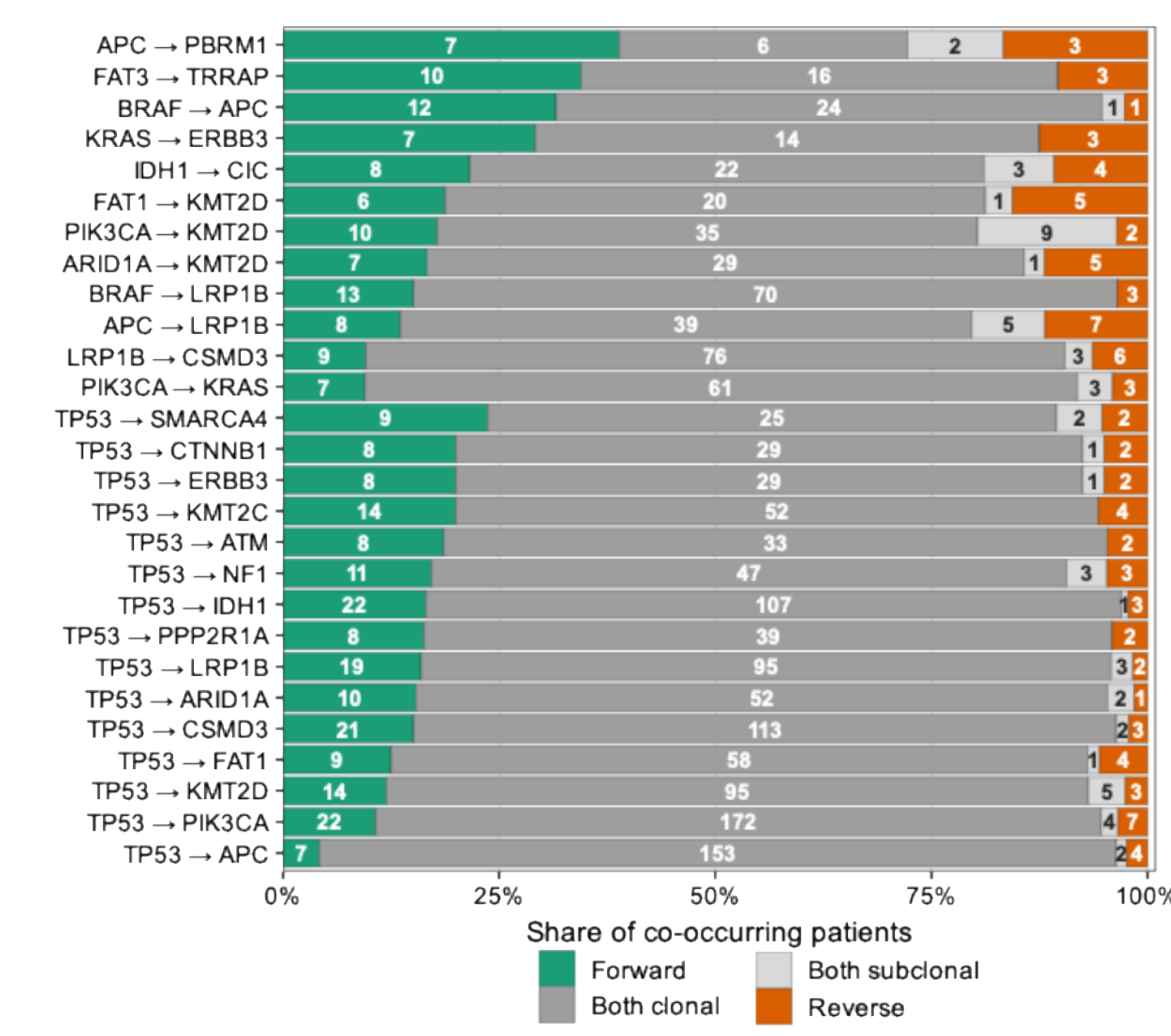
Drivers appeared in 16.7% of subclonal mutation clusters. Gene-level patterns were heterogeneous: TP53, APC, and KRAS were mostly clonal, whereas KMT2C, CTNNB1, and ATM were more often subclonal.

Inferring timing of driver pairs

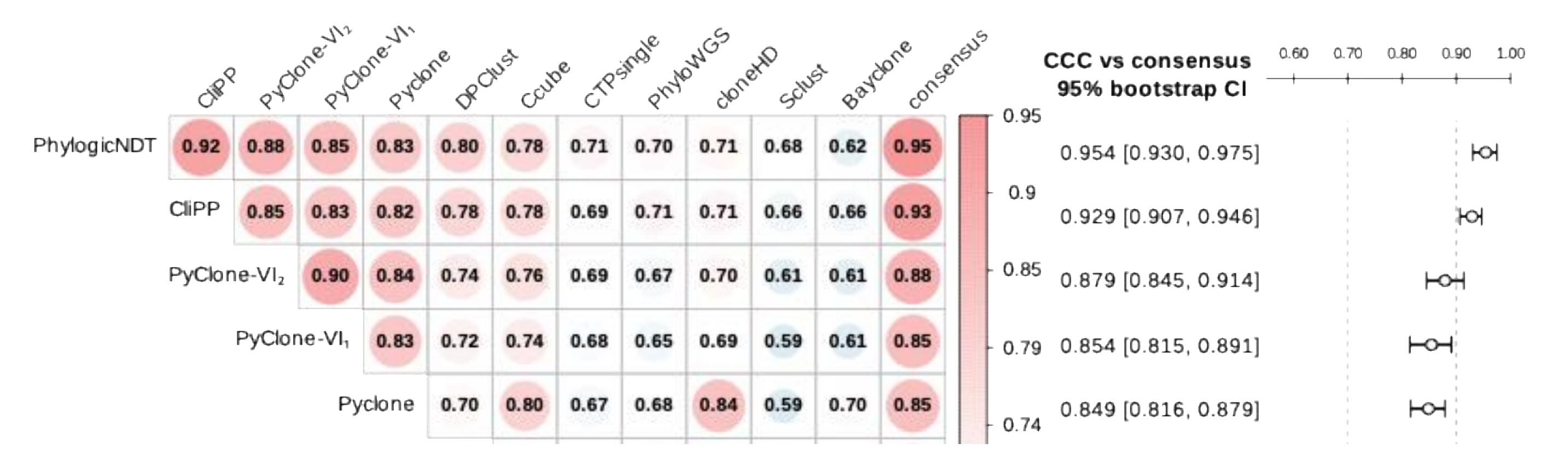
We infer driver order from CliPP clonal/subclonal status: if gene A is clonal and gene B is subclonal in the same tumor, A is interpreted as earlier.

Pairs are classified as Forward, Reverse, Both clonal, or Both subclonal. Plotted pairs have ≥ 10 resolvable patients, where resolvable = Forward + Reverse.

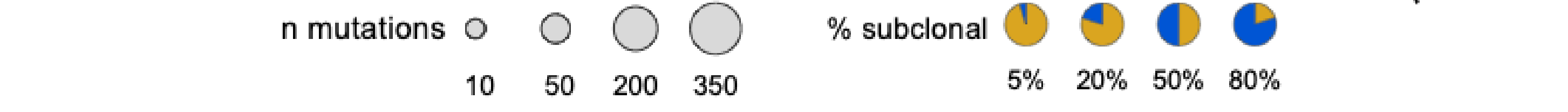
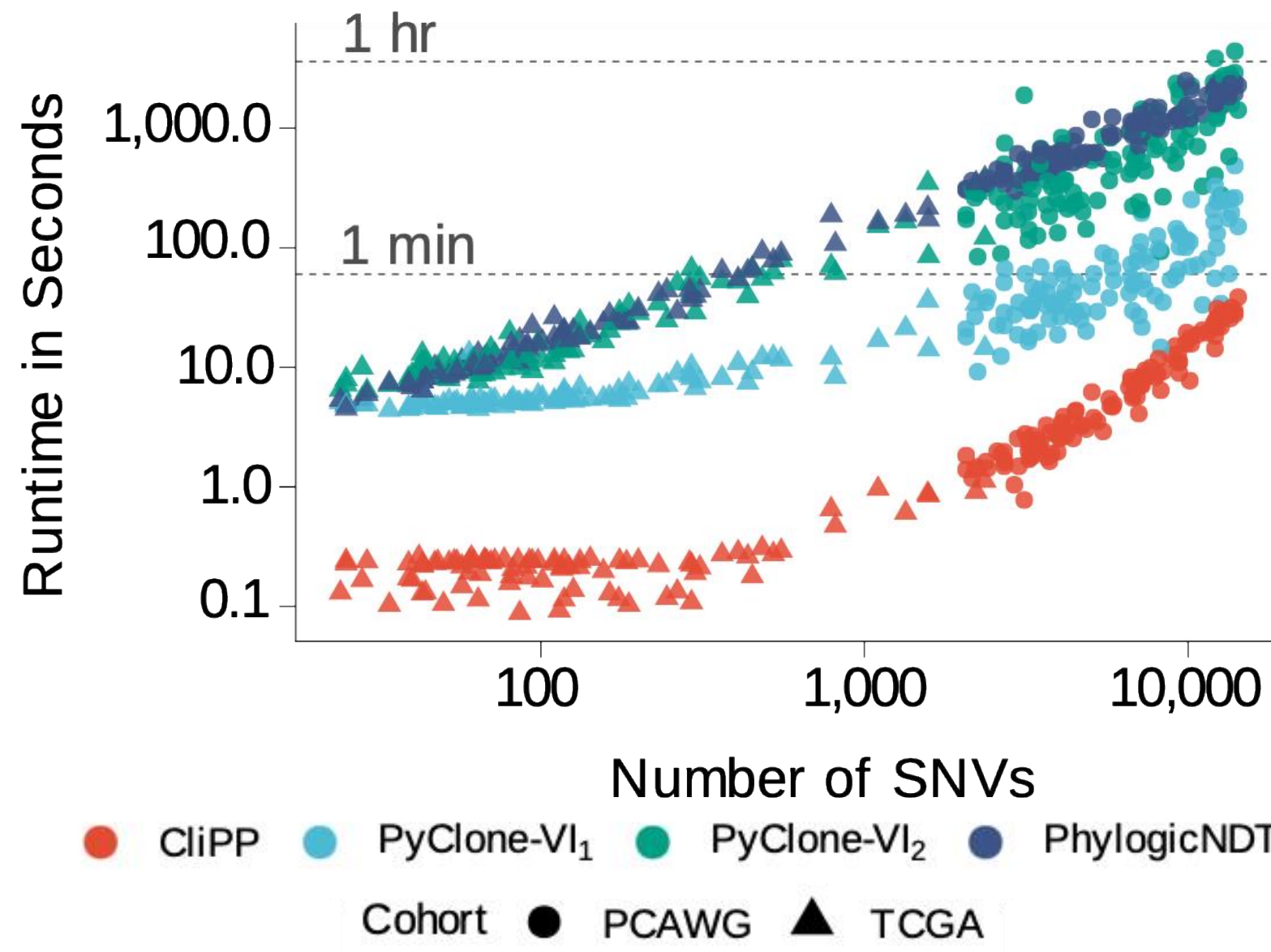
Most pairs were both clonal, reflecting shared early acquisition and limited timing resolution from single-biopsy bulk sequencing, with forward directions capturing known biology.



Benchmarking accuracy in PCAWG data (n=2,008)



PCAWG-consensus	Clonal	6,155 (88.5%)	210 (3.0%)
	Subclonal	52 (0.7%)	539 (7.7%)
	Clonal	PCAWG-CliPP	
	Subclonal		



Each bubble represents one gene x cancer-type pair. Bubble area scales with the number of driver mutations in that cell; pie wedges show the clonal versus subclonal fraction. Cancer types on the x-axis are ordered by median tumor mutational burden, and genes on the y-axis are ordered by total mutation count. The top row summarizes all drivers within each cancer type; the rightmost column summarizes each driver gene across all cancers.

TP53, APC, and KRAS were mostly clonal across cancers, consistent with early founding roles in tumor evolution. In contrast, KMT2C, CTNNB1, and ATM showed higher subclonal fractions, suggesting later acquisition in a subset of tumors.



CliPP-on-web for coding-free subclonal reconstruction

CliPP-on-web provides a coding-free interface for subclonal reconstruction and visualization. Users can upload SNV, copy-number, and purity files; run CliPP; inspect VAF and cellular-prevalence distributions; query TCGA and PCAWG samples; check clonal/subclonal status of TCGA driver mutations; and download CliPP output files.

1. Jiang, Y.J. et al. Scalable subclonal reconstruction of cancer cells in DNA sequencing data using a penalized likelihood model, *bioRxiv*, (2026)
2. Fan, J. & Li, R. Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties. *J. Am. Stat. Assoc.* 96, 1348–1360 (2001).
3. Dentre, S. C. et al. Characterizing genetic intra-tumor heterogeneity across 2,658 human cancer genomes. *Cell* 184, 2239–2254.e39 (2021).
4. Dentre, S. PCAWG Intra-Tumor Heterogeneity Simulations. 1. (2021).
5. Martínez-Jiménez, F. et al. A compendium of mutational cancer driver genes. *Nat. Rev. Cancer* 20, 555–572 (2020).