

From Exome to Whole Genome sequencing analysis

Florian SANDRON^{1,3}, Lilia MESROB^{1,2}, Nicolas WIART¹, Aurélie LEDUC¹, Ghislain SEPTIER¹, Stéphane MESLAGE¹, Delphine BACQ¹, Vincent MEYER¹, François ARTIGUENAVE¹

¹CEA/LBI (Laboratoire de Bioinformatique et Informatique), ²INSERM, ³INRA/France Génomique
 Centre National de Génomique Commissariat à l'énergie Atomique et aux Énergies Alternatives
 Contact: vincent.meyer@cng.fr

1. Varscope: High throughput calibrated process

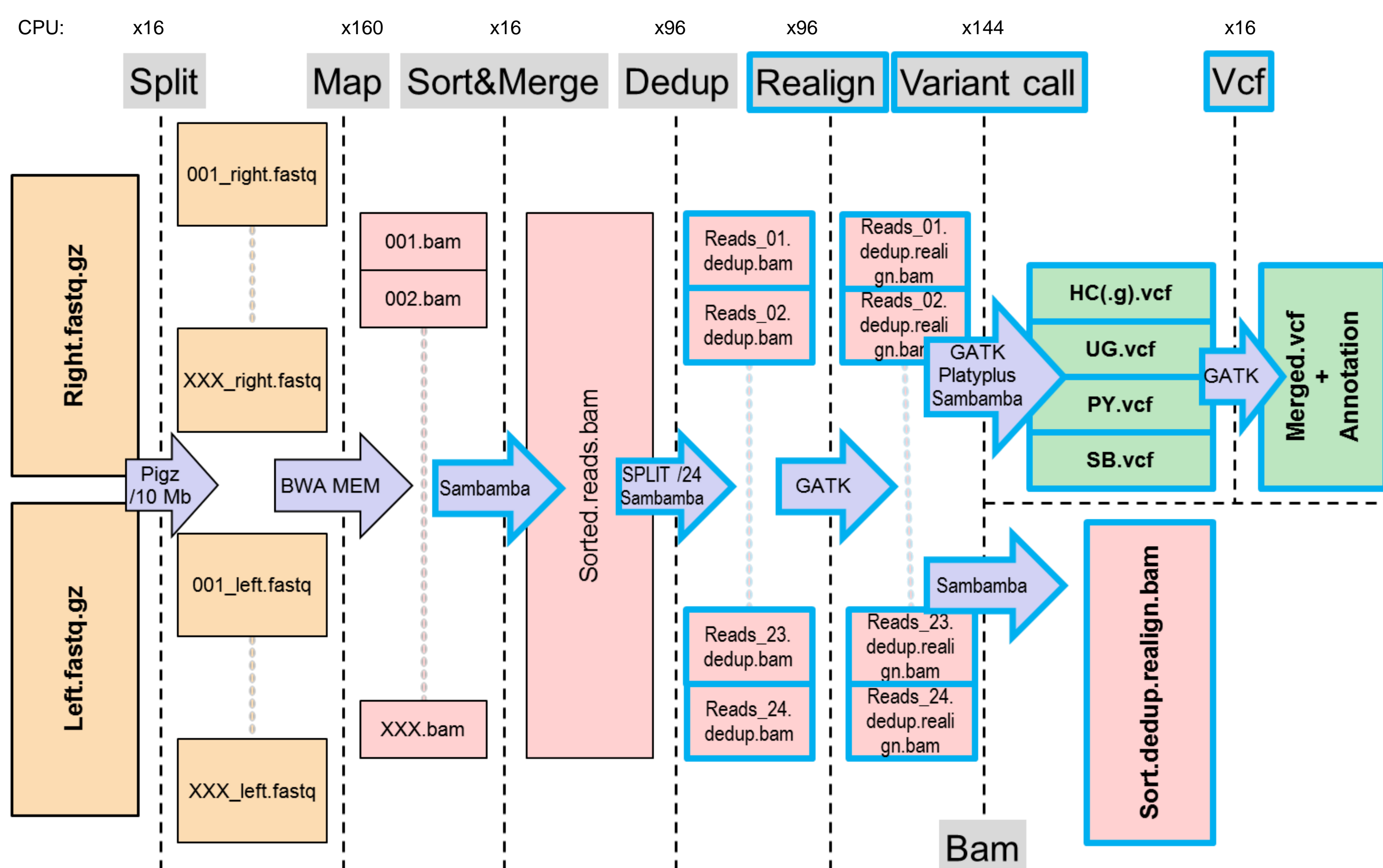


Fig. 1: Overview of the CNG/CCRT mapping and variant calling pipeline organization. Last developments are highlighted in blue

	#base	Fastq.gz	bam	Vcf	#variant
Exome 70x	8,7 Gb	9,7 Go	5,7 Go	16 Mo	75 483
WG 30x	98 Gb	84 Go	75 Go	1075 Mo	4 792 544
WG 70x	223 Gb	196 Go	187 Go	1166 Mo	4 958 931

Fig. 2: Sizes/Values summary related to the files generated during WES and WGS analysis

- Optimization of the « Map-reduce » strategy (fig.1)
- Softwares upgrade and update (fig.1)
- Addition of a « multical » strategy (fig.1,3)

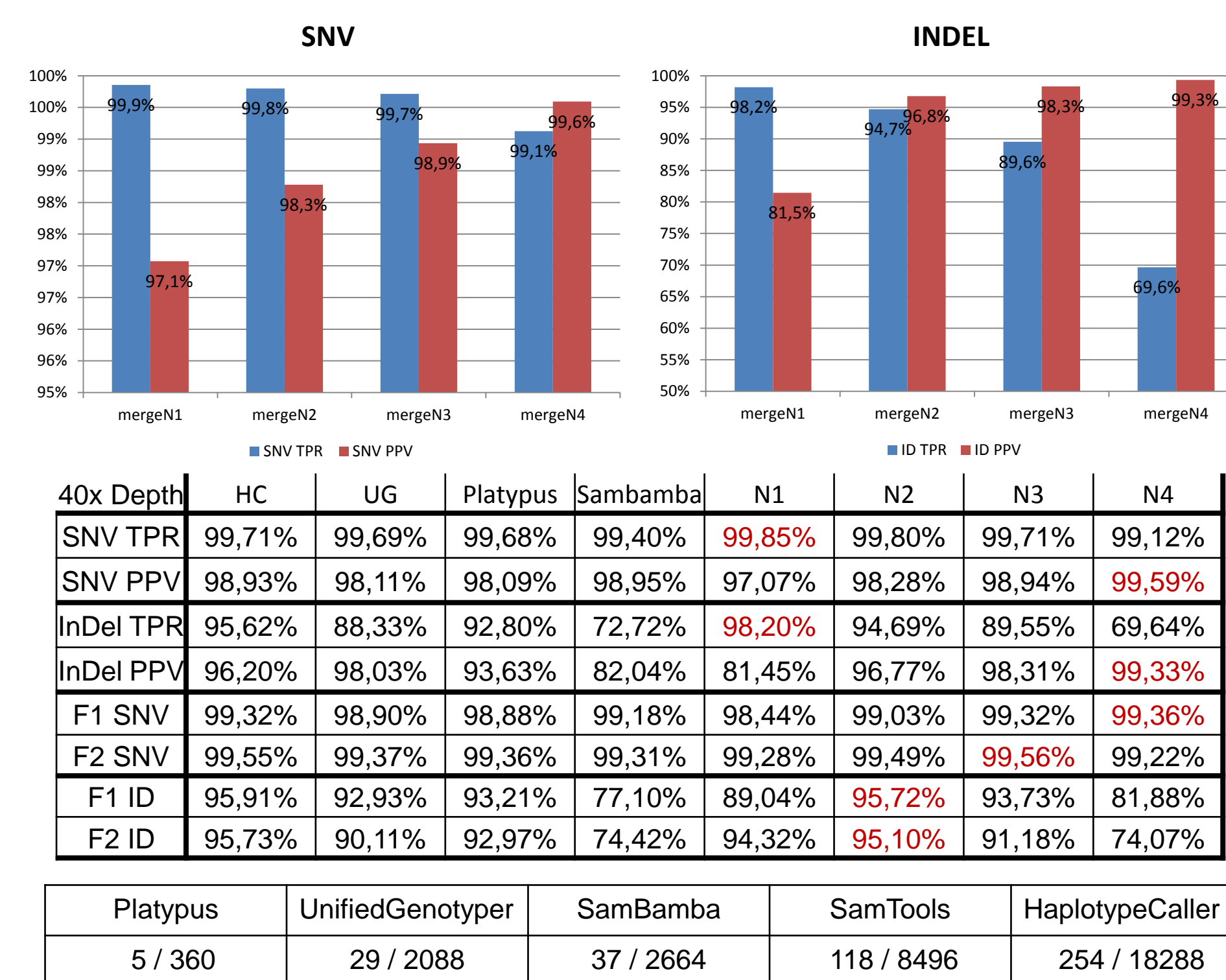


Fig. 3: Above calling program and « multical » strategy performance (NIST/GIAB.v2.19) and approximate Real-Time/Machine-Time(min/72cores)

Abstract

The last few years show an enthusiastic evolution of the human genome resequencing projects profile. We reached the point where Whole Genome Sequencing (WGS) projects are about to become a new standard for research and personal genomics. To achieve this goal, major challenges have to be addressed: Management of massive data storage (1Po/year for X5 illumina), Optimization of computing resources and development of new analysis strategies. Here we present a focus on the last progress carried out at the Centre National de Génomique to handle this leading transition phase.

Since 2013, bioinformatics analyses of the CNG are processed at the HPC facility TGCC: CEA, "Très Grand Centre de Calcul", Bruyère le Chatel, 2 Pflops including 400 Tflops and 5 PBytes dedicated to France Genomics projects. For example, pipelines for exome sequencing data analysis (varscope), mainly based on best practices for mapping, calling and first pass annotation (DePristo et al., 2011/Van der Auwera et al., 2013) was setup in 2013 but required a major upgrade to ensure speed and scalability fitting the production flow anticipation (up to 9000 WG by year for a X5 illumina sequencer).

We present here the high throughput calibrated process deployed on the TGCC facility (fig.1). This process include a "map-reduce" optimization and software upgrade to handle the 10 times increase capacity required by the transition of Whole Exome Sequencing (WES) to WGS (fig.2).

Moreover, we benchmarked several variant calling algorithms (NIST/GIAB; Zook et al., 2014) and added "4-multical" step to our process based on HaplotypeCaller(HC; McKenna et al.,2010), UnifiedGenotyper(UG), Platypus(PY; Rimmer et al., 2014), and Sambamba/Samtools(SB/ST; Tarasov et al., 2015; Li et al., 2009)(fig.1, fig. 3).

We evaluated our process from low coverage analysis (10X, ex: very large cohort sequencing), to standard coverage (30X) and finally high coverage sequencing (> 100x, ex: somatic mutation or mosaicism detection). In "standard" condition, data from a 30X experiment are currently completed in approx. 8 hours (fig. 1). We didn't identify any deadlock situation up to a 240X coverage (equivalent to a full X5 flowcell). Finally, we compared relative performance of WES and WGS and showed equivalent coding variants detection, steadier genome coverage distribution and up to 50% wider access to biologically relevant annotated regions (fig 7,8,9).

In conclusion, we are facing a major transition which requires refactoring and development to harvest knowledge from constantly over-flooding sources of data. We showed here benefits brought by WGS over WES in term of coverage and precision and important upgrades realized to support large WGS programs.

Our next step in this context will be to evaluate and add compression data strategies for storage and processing optimization. In addition we will add our structural variation process to the production flow to provide our collaborators a more complete view of the genome organization.

2. Scalability: from low to high coverage

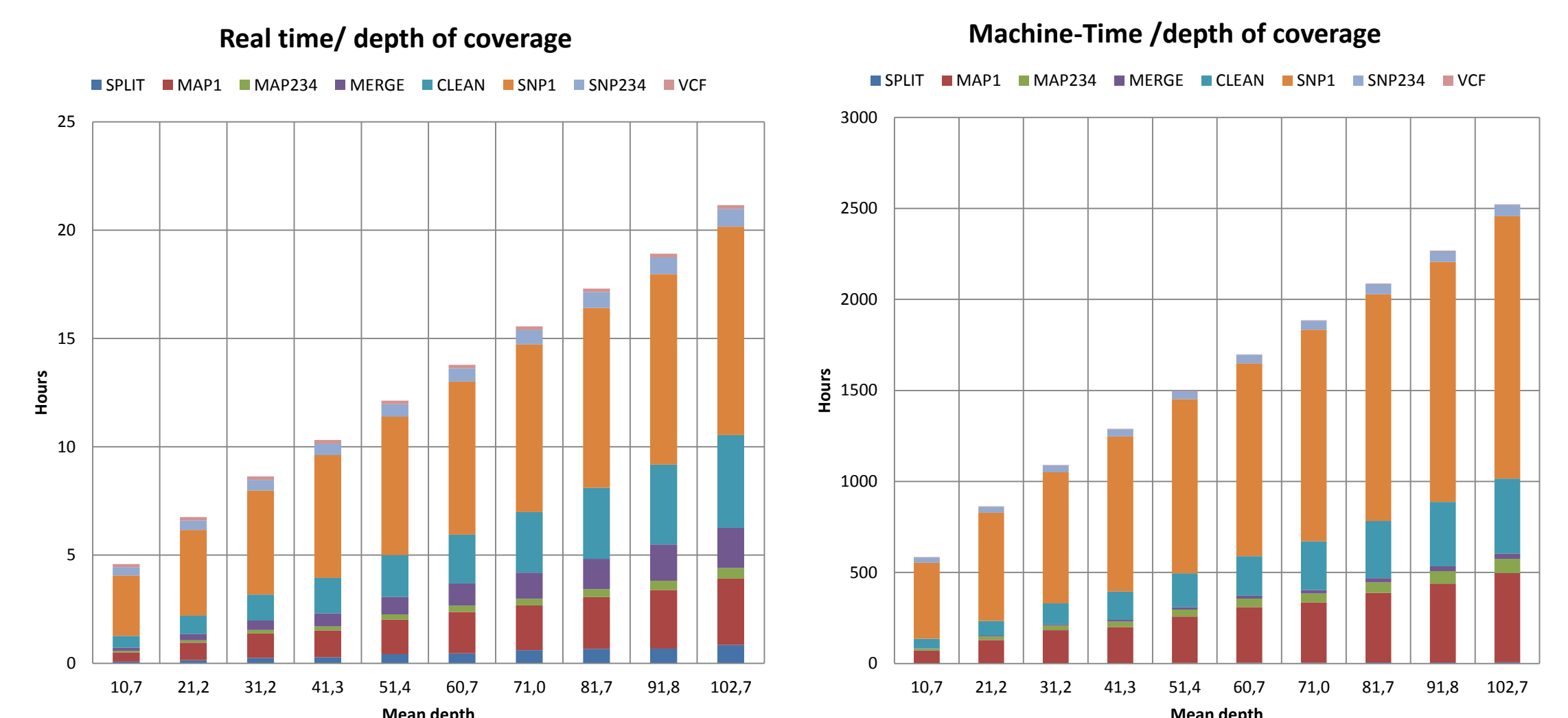


Fig. 4: Computing time for each major analysis step related to the depth of coverage

3. Exome to Whole Genome transition

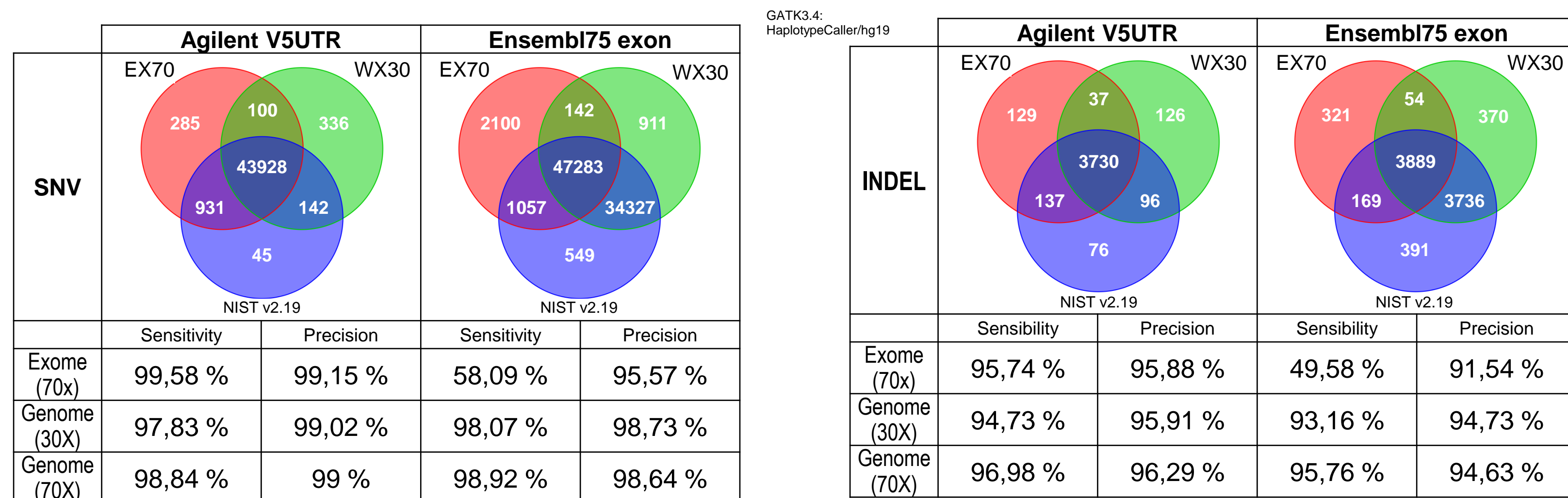


Fig. 7: Variant calling comparison between WES and WGS on Agilent V5UTR exome targeted sequencing and Ensembl « exons » selection(NIST/GIAB)

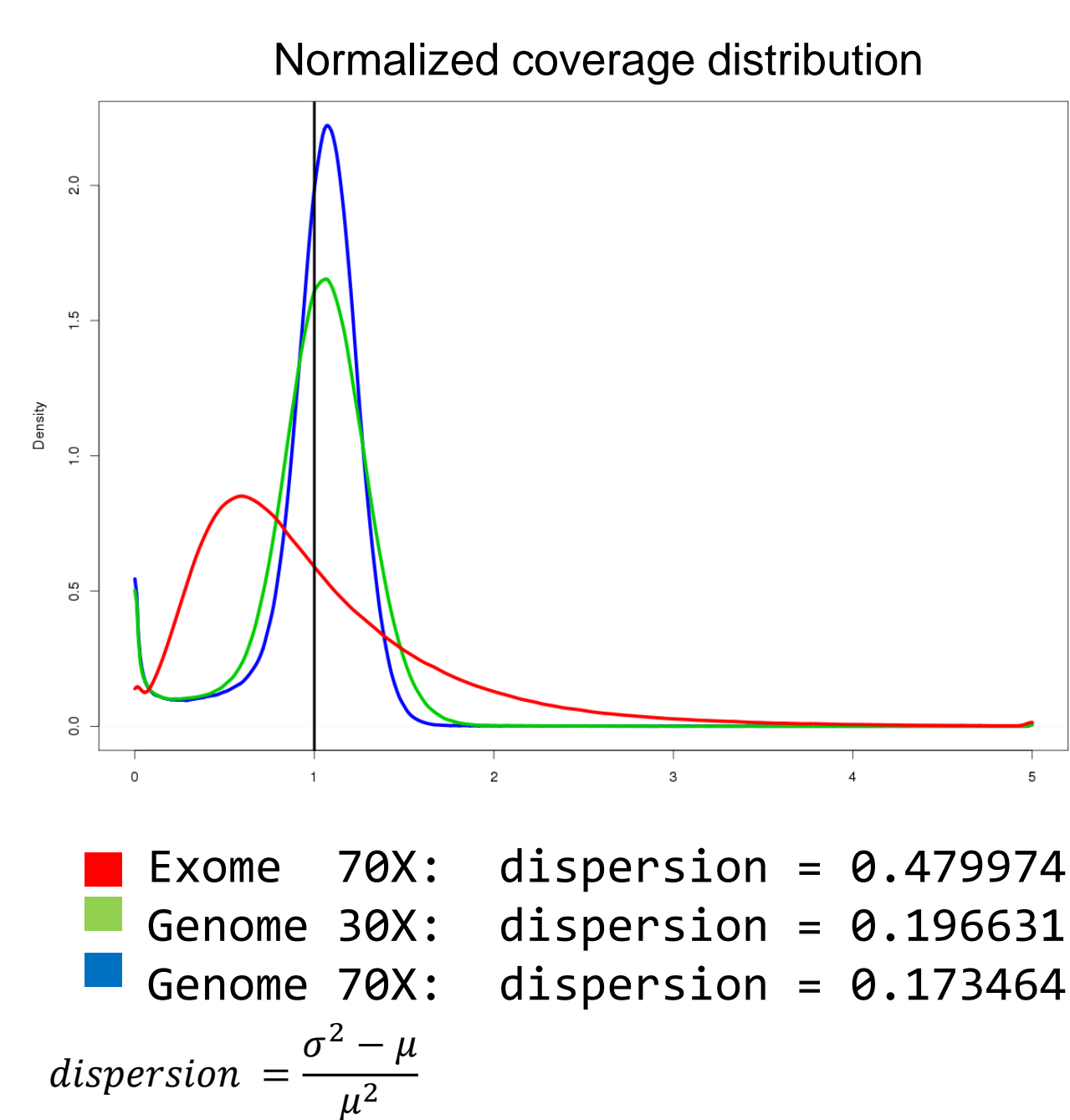


Fig. 8: Comparison of Coverage dispersion between WES and WGS on AgilentV5-UTR targeted regions

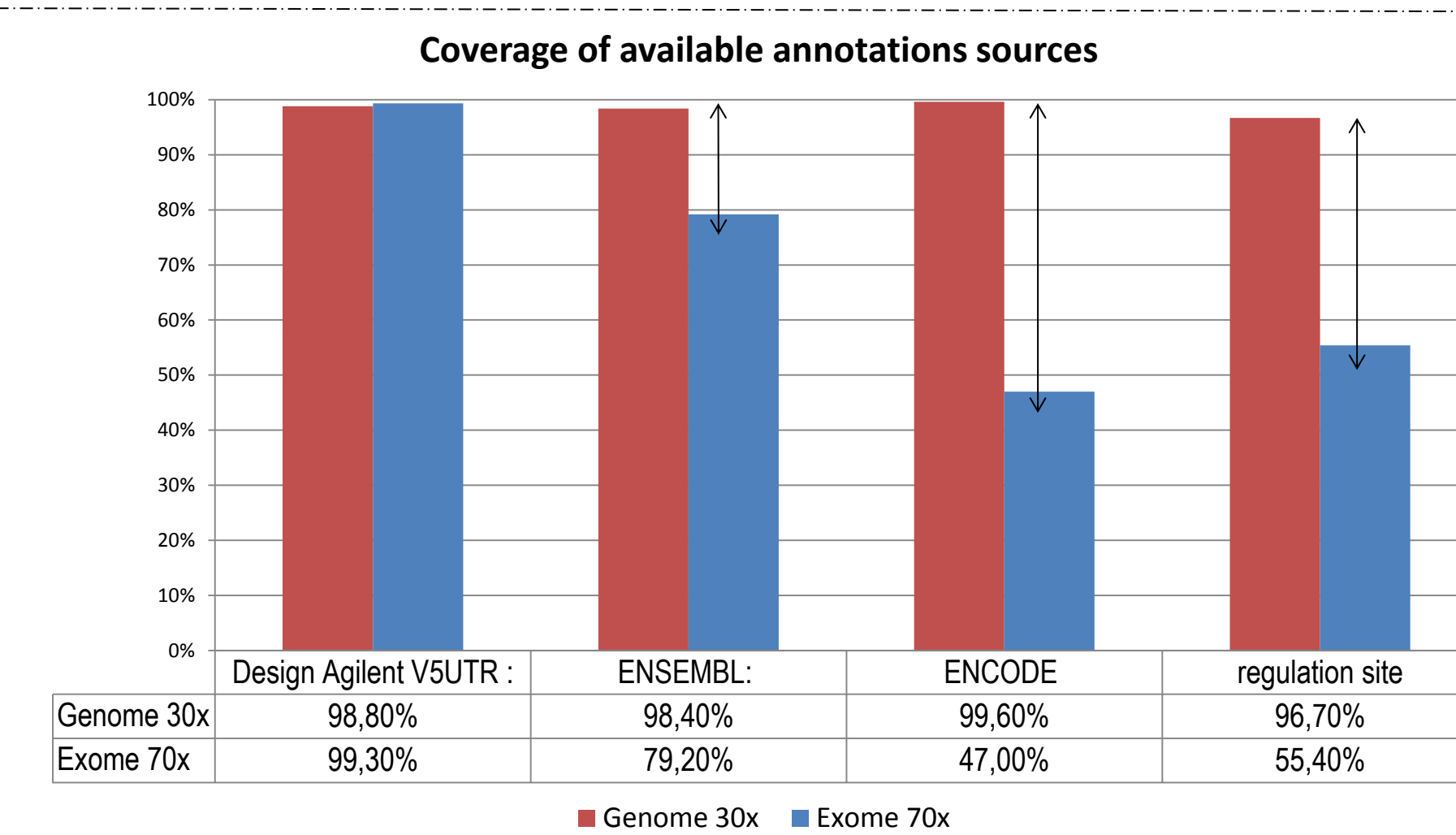


Fig. 9: Annotation coverage comparison between WES and WGS according to Ensembl 75/GRCh37.p13

- WGS Calling performance mostly similar on exome targeted regions (fig.7)
- Steadier coverage for WGS (fig.8)
- Wider annotation access (fig.7,9)

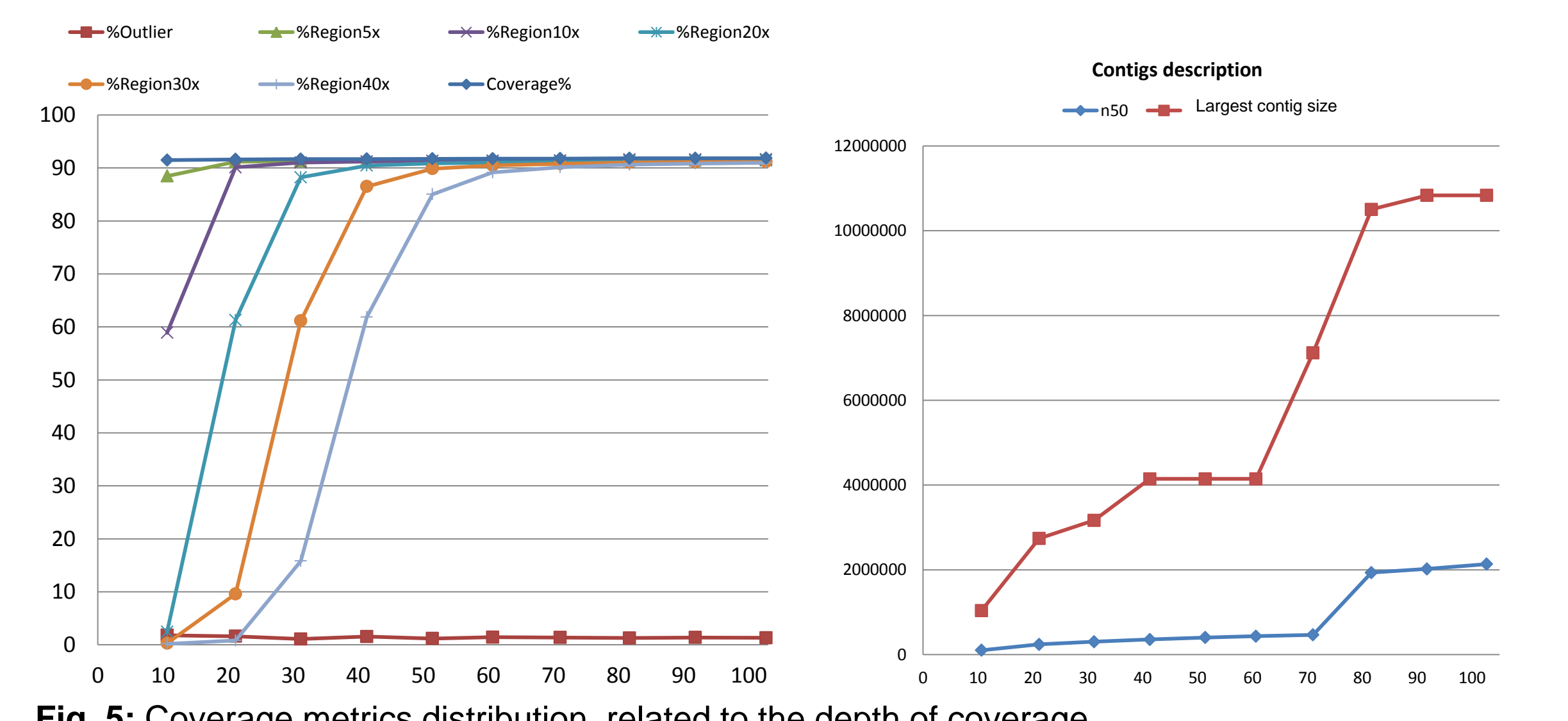


Fig. 5: Coverage metrics distribution related to the depth of coverage

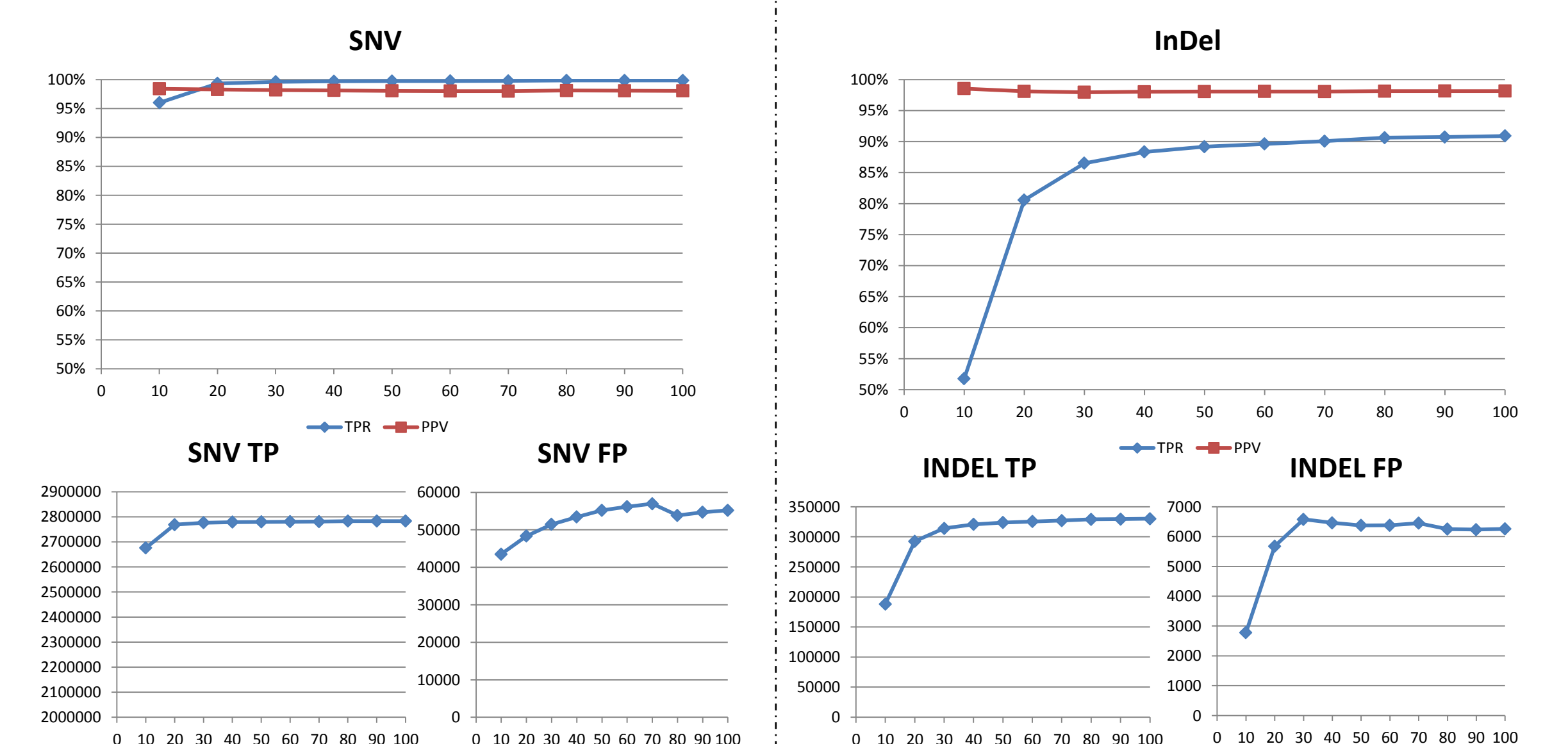


Fig. 6: Evolution of the indel/SNV sensitivity/precision related to the depth of coverage (NIST/GIABv2.19)

- 2,5h for a WES 60X
- 8h for a WGS 30x up to 42h for a WGS 240x
- 100 WES (70X) / ~12 WGS (30X) overnight

Perspectives

- Addition of compression data strategies (cf. JY.Lalanne; N.Wiart)
- Switch to GRCh38 as default analysis environment
- Addition of our structural variation detection process to the production flow