# CNV detection

# from targeted next-generation sequencing data

**Anna Benet-Pagès**

**Johor 29.08.2018**

# Detection of structural variants and human disease
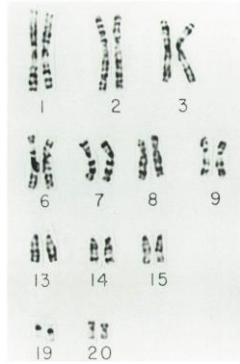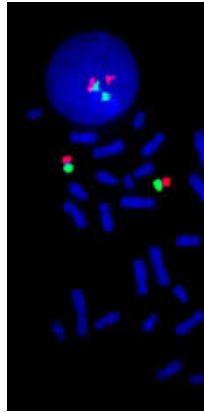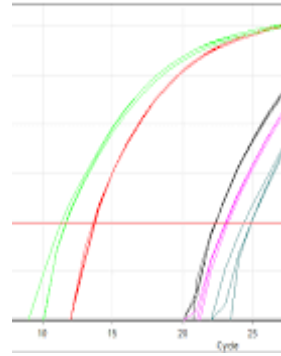


Calvin Bridges, 1936
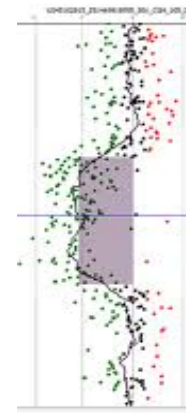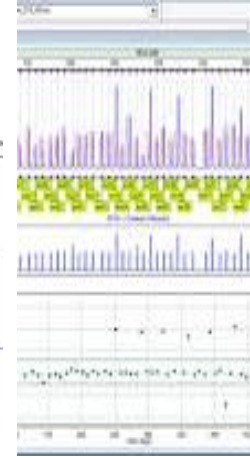*B/B+*

**G-banding** Rowley et al., 1973
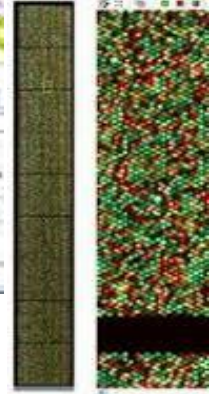
**FISH** Langer et al., 1982

**Quantitative PCR** Porcher et al., 1992

**Array CGH** Pinkel et al., 1998c

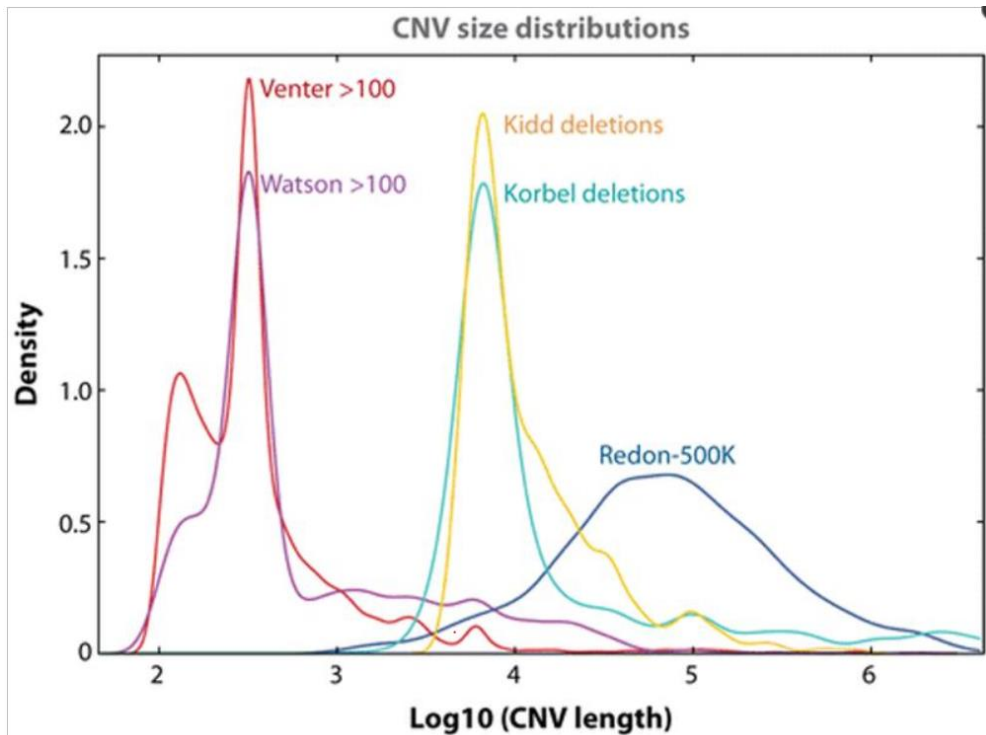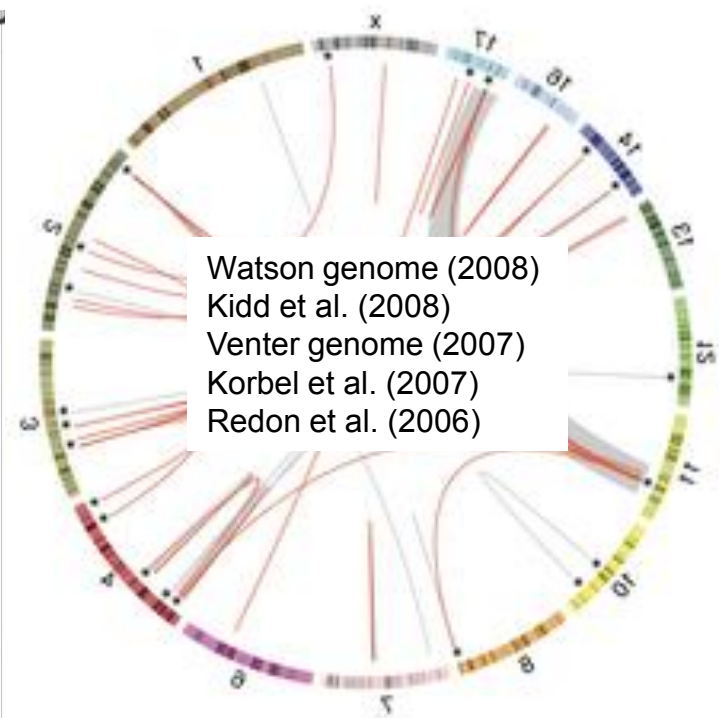**MLPA** Schouten et al., 2002

**SNP array** Komura et al., 2006

# High resolution technologies reveal small-size CNVs

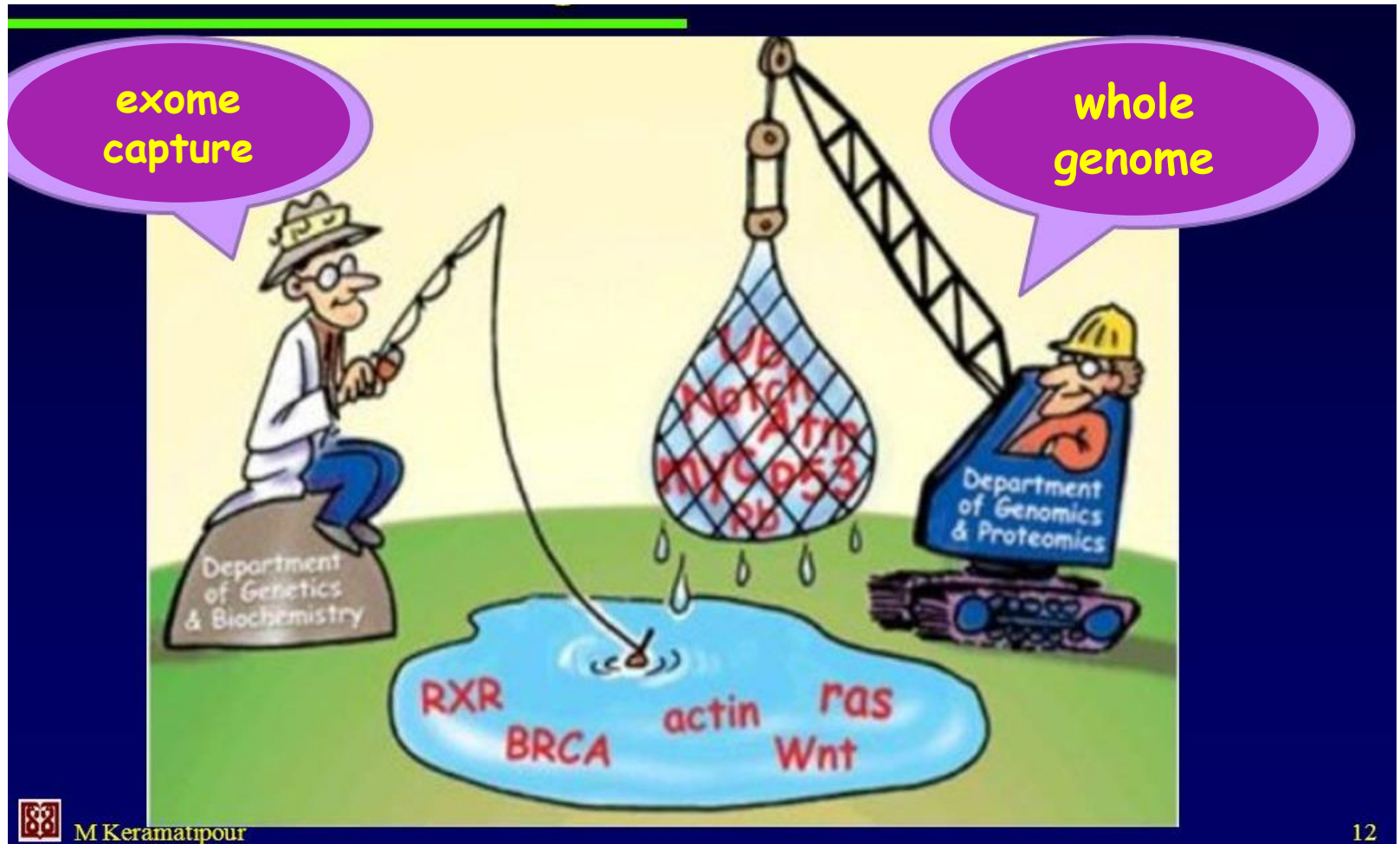Size distribution of copy number variations (CNVs) larger than 100 bp



*Zhang et al. 2009*

Watson genome (2008)
Kidd et al. (2008)
Venter genome (2007)
Korbel et al. (2007)
Redon et al. (2006)

➤ Smaller structural variants are the most frequent

# On the meanwhile…

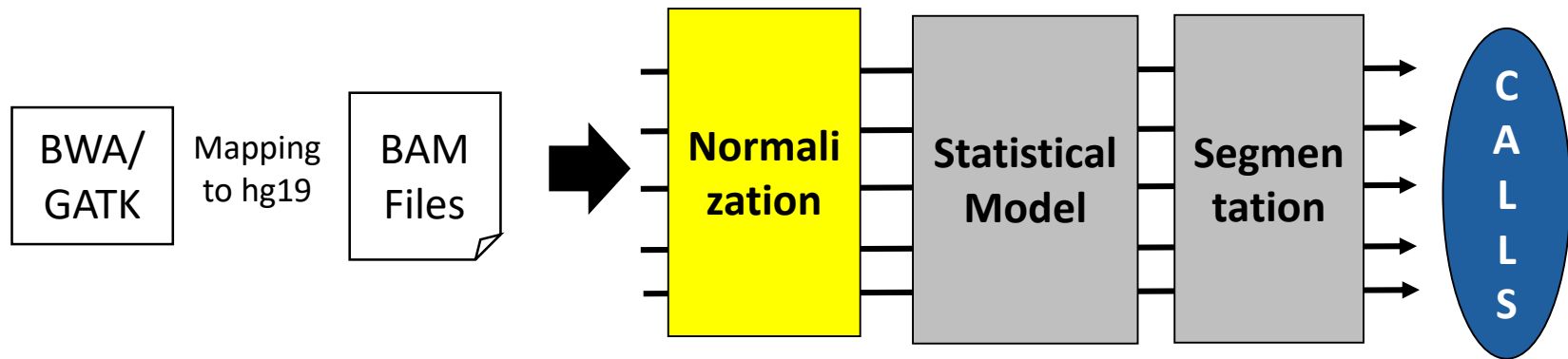# CNV detection from targeted-capture data

## Challenges:

- CNV detection from exon capture approaches depends solely on read depth data

- Enrichment efficiency introduces a systematic noise in read depth data

- Coverage bias between sequencing runs and within samples of the same run

- Single exon events are extremely difficult to detect

- Control individuals are difficult to obtain (reference set / validation)

- Validation is expensive

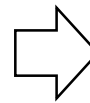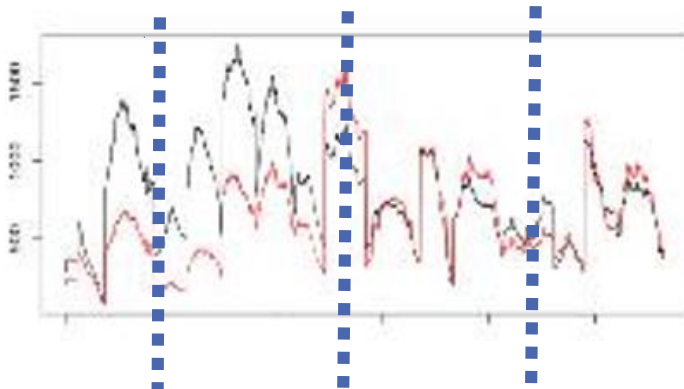# CNV detection methods general considerations

➢ **Which tool should I choose?**

- applicable to capture data

- calling of rare CNVs

- easy to integrate (take bam files as input)

- easy handling (installation / running time)

- multi-sample usage (possibility to normalize against reference set)

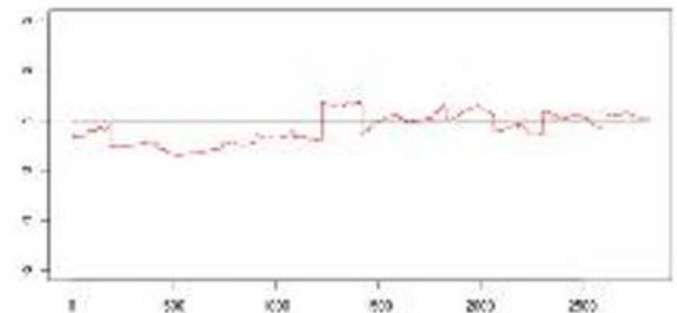- Tools should use different statistic models

# CNV Pipeline Structure



```
BWA/          Mapping      BAM
GATK          to hg19      Files
```

**Normalization** → **Statistical Model** → **Segmentation** → **CALLS**

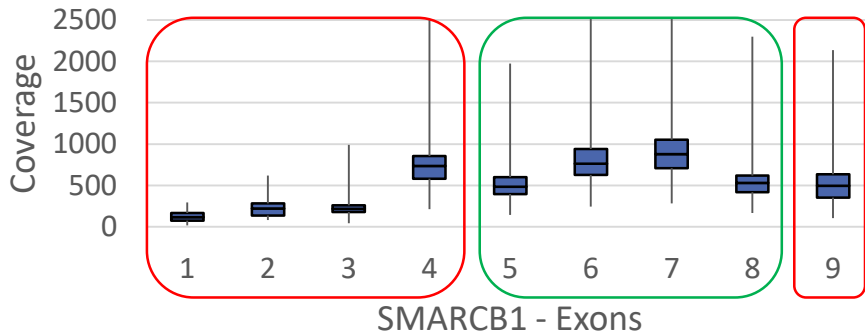breakdown of the target region exons/ windows

GC
Mapping qual
# reads
Ref. Set

# Reference Sets and data normalization

➢ different reference sets for different kits / enrichment methods

➢ normalization against samples from the same sequencing run to improve robustness against workflow conditions
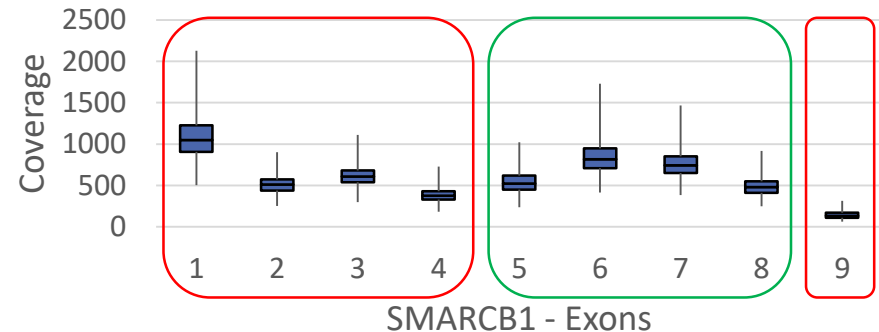
# CNV Pipeline Structure
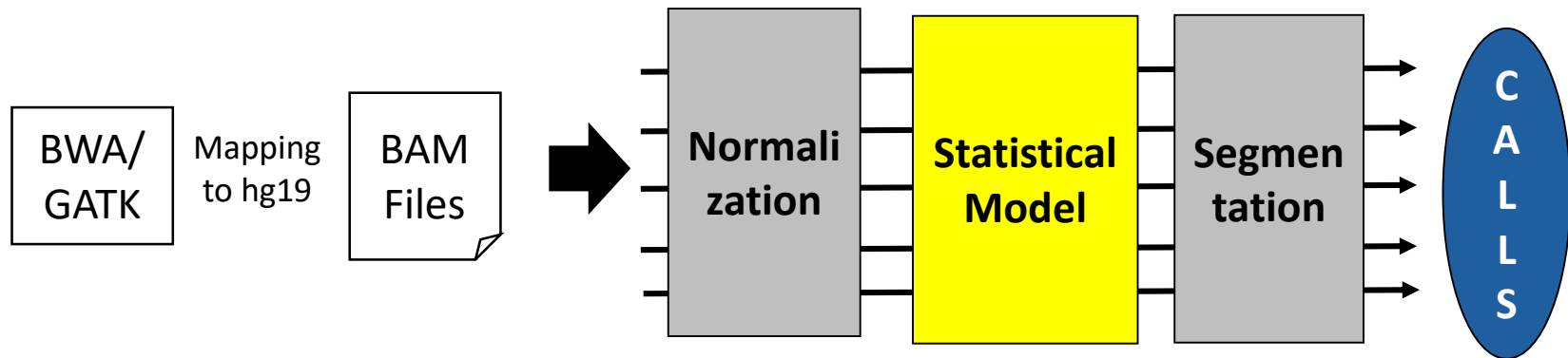


| | Mapping to hg19 | | | Normalization | Statistical Model | Segmentation | CALLS |
|---|---|---|---|---|---|---|---|
| BWA/ GATK | | BAM Files | | | | | |

Poison distribution
Beta binomial
Negative binomial
Normal distribution

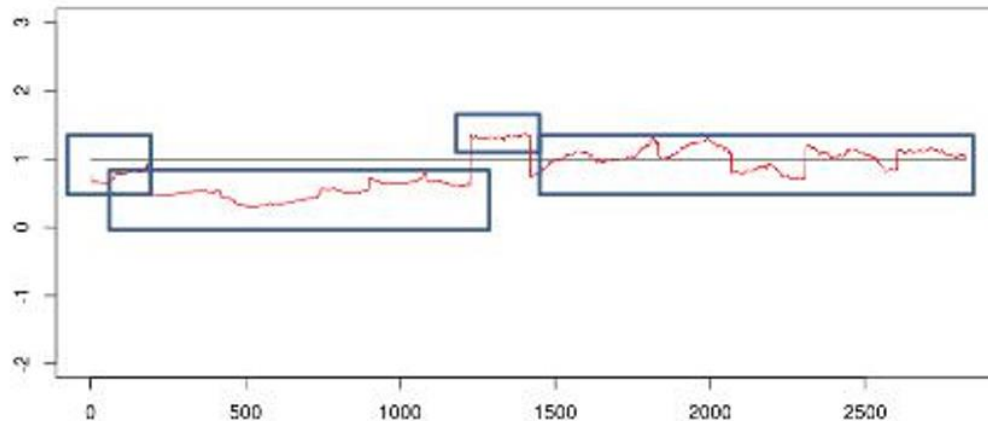| Coverage | frequency |
|---|---|
| 50 | 24 |
| 100 | 95 |
| 200 | 82 |
| 500 | 21 |

# CNV Pipeline Structure



grouping areas (exon/window) with the same prediction (gain / loss / normal)

# CNV detection methods

> Use a combination of several detection tools

AGE, BicSeq, BreakDancer, Breakpointer, Breakseq, Canoes, Clamms, Clever, ClipCrop,

Cn.MOPS, CNAnorm, CNAseg, CND, CNV_TV, Cnvator, CNVer, CNVer, HugeSEQ,

hydra, inGAP_sv, JointS...    **„meta-CNV-caller"**    rcanavar, Patchwork, pemer

,ReadDepth, rSW_seq, segseq, seqcbs, CNVer, cnvHiTSeq, cnvrd, CNV-seq,

conserting, CONTROL_FREEC, cops, copySeq, crest, ERDS, codex

EWT_RDXplorer, GasvPRO, GENSENG, XHMM

*Noll et al., Npjgenmed 2016*

# Meta-Tool CNV Detection Pipeline

- **ExomeDepth**

  extremely sensitive and robust against samples that do not correlate with the reference

- **Canoes**

  has a high sensitivity for small deletions, high performance in low coverage regions and with few reference samples

- **Clamms**

  corrects for GC content and mappability, divides large exons into smaller regions and calls also common CNVs

- **Codex**

  corrects for GC content and mappability, calls also common CNVs, uses no HMM for segmentation (all other tools use HMMs)

- **Inhouse method**

  is well adapted on inhouse data, screens for heterozygosity, corrects for GC content, exon score depends on previous analyses

# Performance of single tools

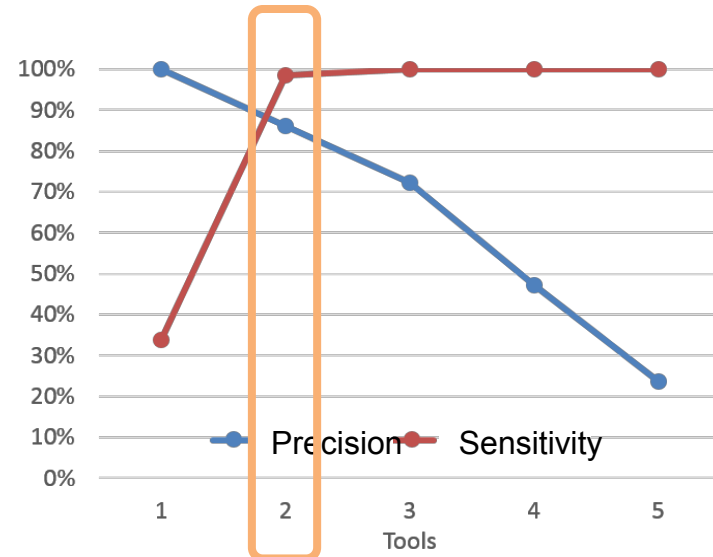➢ Training set: true set of 146 CNV calls detected via MLPA.

|  | Exome Depth | Clamms | Canoes | Codex | In-house |
|---|---|---|---|---|---|
| **Precision** | 45.63% | 68.57% | 96.77% | 64.75% | 40.82% |
| **Sensitivity** | 90.38% | 46.15% | 57.69% | 63.46% | 76.92% |

# Performance of tool combinations

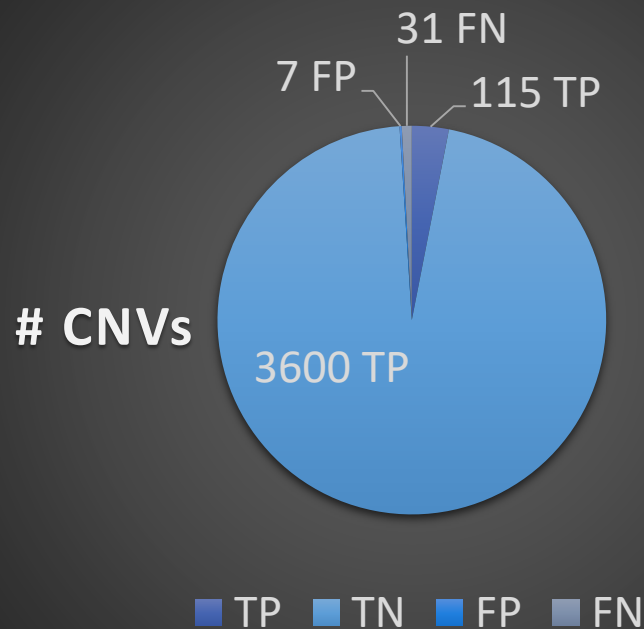➢ What is the best number of tools required to call a variant?

| stringency | | TP | FN | FP | TN | Sensitivity (TPR) | Specificity (TNR) | Precision (PPV) | NPV |
|---|---|---|---|---|---|---|---|---|---|
| | **2 out of 5** | 115 | 7 | 8 | 3600 | 94.26% | 99.78% | 93.50% | 99.81% |
| | **3 out of 5** | 72 | 50 | 1 | 3600 | 59.02% | 99.97% | 98.63% | 98.63% |
| | **4 out of 5** | 28 | 94 | 1 | 3600 | 22.95% | 99.97% | 96.55% | 97.46% |
| | **5 out of 5** | 4 | 118 | 0 | 3600 | 3.28% | 100.00% | 100.00% | 96.83% |

➢ Using two out of five concordant tool predictions shows the best balance between sensitivity and specificity

# CNV Pipeline Evaluation

- >3700 MLPAs were performed in ~90 genes
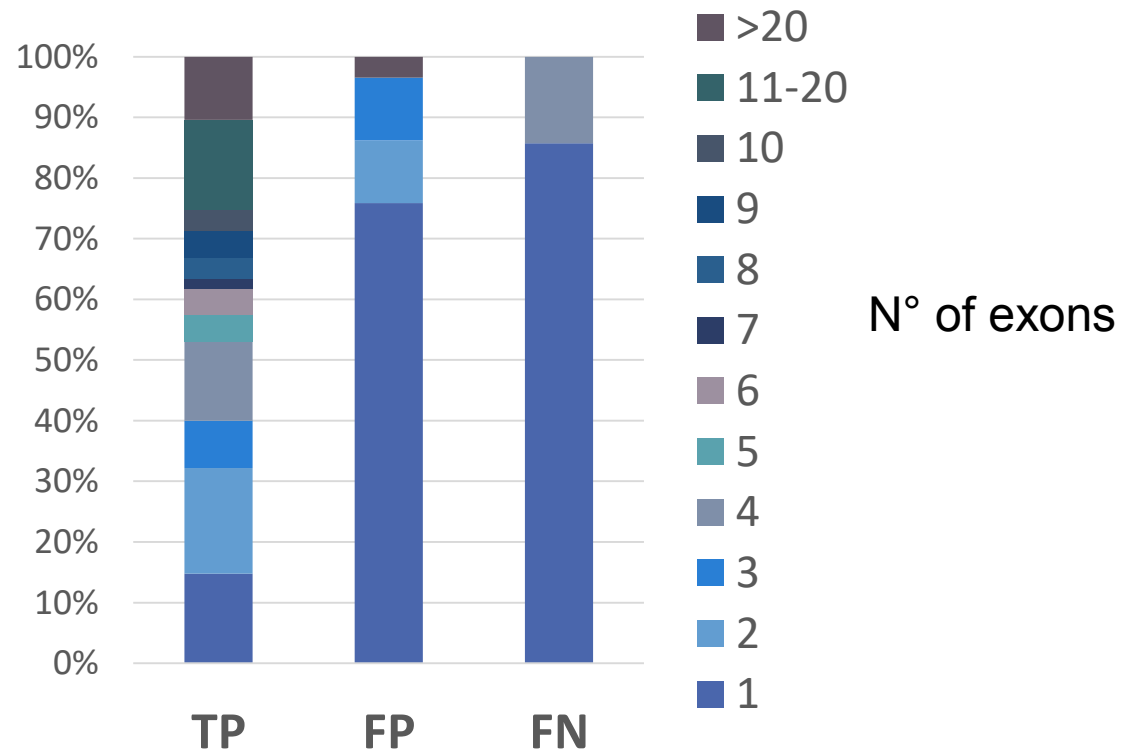- 146 CNVs (85 deletions / 61 gains)



**# CNVs**

- 31 FN
- 7 FP
- 115 TP
- 3600 TP
- ■ TP  ■ TN  ■ FP  ■ FN

Sensitivity:      88.60%
Specificity:      98.88%
Precision:        71.40%

**Pseudogenes excluded:**

Sensitivity:      94.26%
Specificity:      99.78%
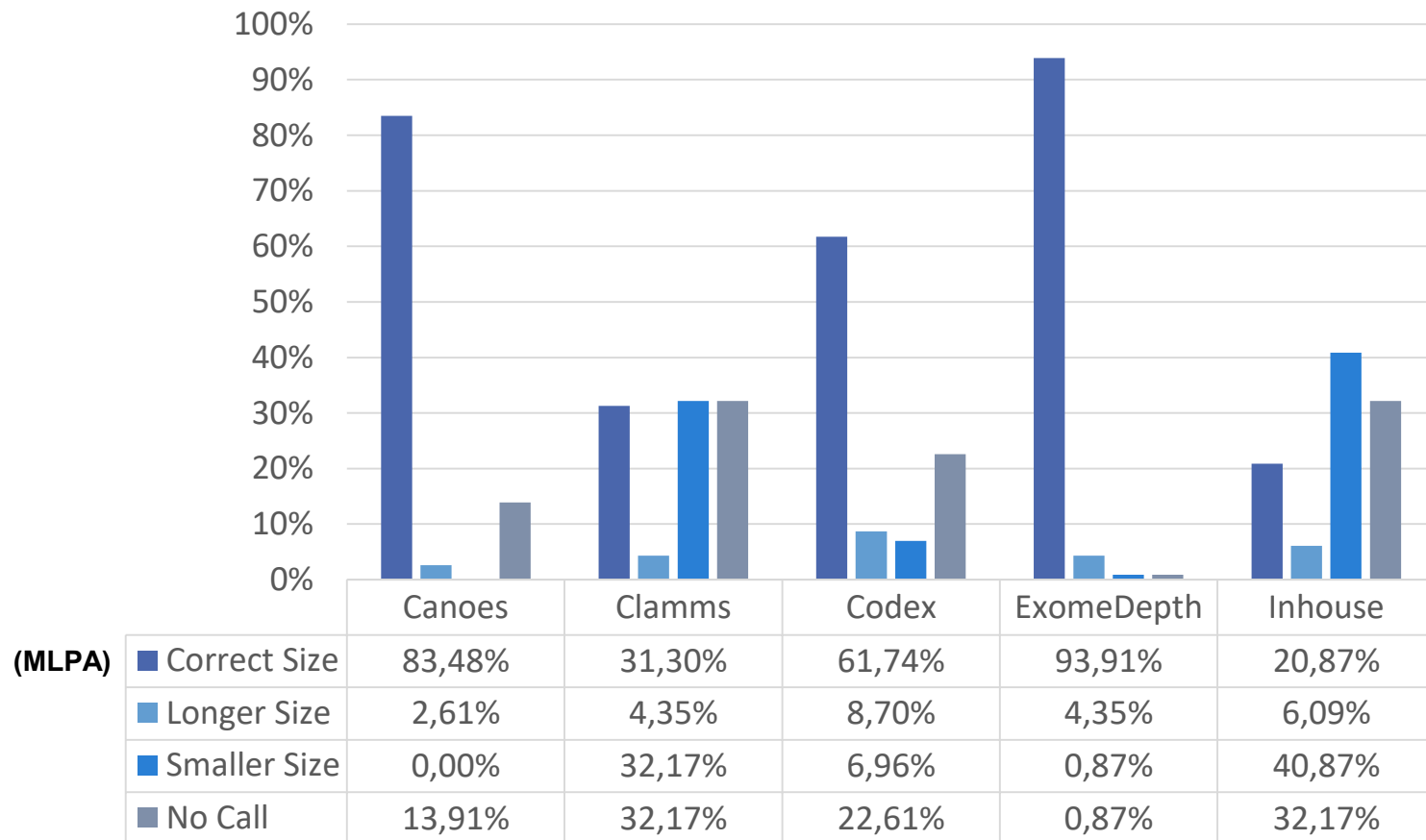Precision:        93.50%

# TP, FP , FN versus CNV size

➢ Comparison of CNV sizes of TP, FP and FN calls detected by the combined CNV pipeline on the validation set.

➢ FP mainly cosiand FN calls consist mainly of single exon events.



**definition of special quality thresholds for single exon events to minimize false negatives**

# Discrepancies in CNV size detection between tools

➢ Size of CNV calls were compared to MLPA

➢ Size is given as the number of exons within the call

| (MLPA) | Canoes | Clamms | Codex | ExomeDepth | Inhouse |
|---|---|---|---|---|---|
| ■ Correct Size | 83,48% | 31,30% | 61,74% | 93,91% | 20,87% |
| ■ Longer Size | 2,61% | 4,35% | 8,70% | 4,35% | 6,09% |
| ■ Smaller Size | 0,00% | 32,17% | 6,96% | 0,87% | 40,87% |
| ■ No Call | 13,91% | 32,17% | 22,61% | 0,87% | 32,17% |

# Meta-CNV-caller: multi calls for one event



## Copy Number Variants in NBN

Calls  2

| | Exons | Type | CN | Sample | Pool | Region | PPL | Overlap (min) | Overlap (equal) | Overlap (ExAC) | Overlap (Pool) | Methods |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | E2 – E12 | + | 3 | 121258 | SP-666 | chr8: 90,955,481 - 90,996,789 | | 1 | | 9 | | clamms CN3, exomedepth CN3 |
| 2 | E3 – E15 | + | 3 | 121258 | SP-666 | chr8: 90,947,810 - 90,995,083 | | 1 | | 9 | | canoes CN3, exomedepth CN3 |

# Challenges

❖ Non-uniform of coverage

❖ CNV calling in homologous genomic regions (pseudogenes...)

❖ Clinical interpretation

# non-uniform coverage = capture bias

➢ identification of reliable regions by assessment of capture efficiency to minimize false positives



can not be analyzed

# CNV calling artifacts



2. SCN1A

| | |
|---|---|
| **Number of exons** | 26 |
| **Duplications in** | 5 exons and 4 calls (frequency = 0.00) |
| **Deletions in** | 26 exons and 39 calls (frequency = 0.01) |

Depth of coverage uniformity reference set
Depth of coverage uniformity sample

# dup calls
# del calls

Calls / exon

# CNV calling artifacts

# CNV calling in Pseudogenes



Legend:
- Forward read / unique mapping
- Reverse read / unique mapping
- Non-unique mapping

1. PMS2

exon 15    14   13     12     11

[0 - 789]

[0 - 941]

PMS2 exons 11 – 15 can not be analyzed

# CNV calling in Pseudogenes



PMS2

PMS2CL

# How to identify regions affected by pseudogenes

> ➢ alignments of the human genome with itself using blastz



**Human Chained Self Alignments**

# Interpretation of CNV calls – population DB

➢ Deletions and duplications called based on read depth using XHMM; *Fromer et al.*
➢ Z score for the deviation of observed counts from the expected number



Positive Z scores indicate that the gene had fewer variants than expected.
Negative Z scores are given to genes that had a more variants than expected.

# Interpretation of CNV calls – clinical DB

- ➢ DGV
- ➢ DECIPHER
- ➢ ClinVar
- ➢ ClinGen

# Interpretation of CNV calls – clinical DB

# CNV analysis on **1600** individuals within the routine Dx

CNV clarified the underlying phenotype in 8 % of the cases

**Increase of the diagnostic yield in 3%**

Pie chart labels:
- XLR 1
- XD/XR 2
- XLD 1
- AR 9
- AD 16
- AD/AR 11

benet-pages@mgz-muenchen.de