

# The Sentieon Genomic Tools Improved Best Practices Pipelines for Analysis of Germline and Tumor-Normal Samples

Andreas Scherer, Ph.D. President and CEO

Dr. Donald Freed, Bioinformatics Scientist, Sentieon

**CIOReview**

20 most promising  
Biotech Technology  
Providers

**pharma**  
TECH OUTLOOK

Top 10 Analytics  
Solution Providers

**Gartner.**

Hype Cycle for  
Life sciences



**Golden Helix is a global  
bioinformatics company  
founded in 1998.**



### **Variant Calling**

Filtering and Annotation  
Clinical Reports  
CNV Analysis  
Pipeline: Run Workflows

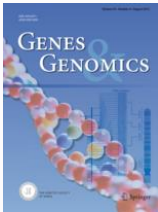
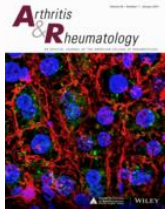
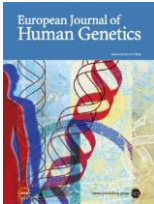
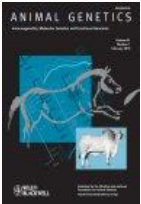


Variant Warehouse  
Centralized Annotations  
Hosted Reports  
Sharing and Integration

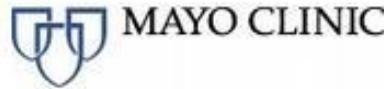


GWAS  
Genomic Prediction  
Large-N-Population Studies  
RNA-Seq  
Large-N CNV-Analysis

# Cited in over 1,100 peer-reviewed publications



Over 350 customers globally

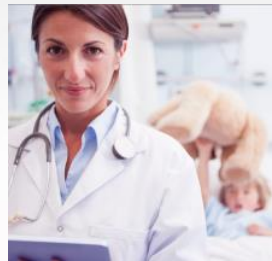


# Golden Helix – Who We Are



When you choose a Golden Helix solution, you get more than just software

- REPUTATION
- TRUST
- EXPERIENCE



- INDUSTRY FOCUS
- THOUGHT LEADERSHIP
- COMMUNITY

- TRAINING
- SUPPORT
- RESPONSIVENESS



- TRANSPARENCY
- INNOVATION and SPEED
- CUSTOMIZATIONS

# End-to-End Architecture for Clinical Testing Labs



**GOLDEN HELIX**  
Enabling Precision Medicine

GENE PANEL    EXOME    GENOME

SEQUENCER

PRODUCTS

BIOINFORMATICS PIPELINE

FUNCTION

 VS-CNV  
 SENTIEON DNASEQ  
 SENTIEON TNSEQ

OMIM    SIFT & POLYPHEN    CLINVAR    ENSEMBL GENES  
CADD    EXAC & GNOMAD EXOMES    DBSNP    REFSEQ GENES  
ONCO MD    CONSERVATION SCORES    COSMIC

FASTQ

SINGLE NUCLEOTIDE VARIATION

BAM

COPY NUMBER VARIATION & LOSS OF HETEROZYGOSITY

VCF

CHROMOSOMAL ABERRATION

ANNOTATED VCF

PUBLIC & COMMERCIAL ANNOTATIONS  
TO ENRICH GENOMIC DATA SETS

CLINICAL REPORT

 **VARSEQ**  
 VSREPORTS  
VSPipeline

ANNOTATE & FILTER  
VISUALLY INSPECT ALIGNMENTS  
VARIANT PRIORITIZATION  
CLINICAL ASSESSMENT

DATA WAREHOUSING

 **WAREHOUSE**

CLINICAL ASSESSMENT CATALOG  
ADVANCED DATA QUERYING  
VERSIONING

WEB-ENABLED INTERFACE  
+ POWERFUL API: JSON, XML  
TSV, CSV, SQL, FHIR

INTEROPERABILITY  
COMPLIANCE WITH HIPAA, CLIA, & CAP  
DATA DISCOVERY



# The Sentieon Genomic Tools – Enabling Precision Data for Precision Medicine



[WWW.SENTIEON.COM](http://WWW.SENTIEON.COM)

# What is Precision Data

---

- Hottest Word: Big Data
- But, PRECISION is the goal:  
precision recommendation, precision prediction
- Big genomics data for precision medicine:
  - precision diagnostics
  - improve precision treatment for individual

Big + Accurate | “better data in, better results out”



**Precision Data**



## Sentieon's Mission

Enable Precision Genomics Data for Precision Medicine

- Ability to process big data at affordable cost and time
- With confidence
  - The highest accuracy
  - Consistent results

# Three components of analytical software

-mathematical methods

-Same mathematical models as the Broad Institute

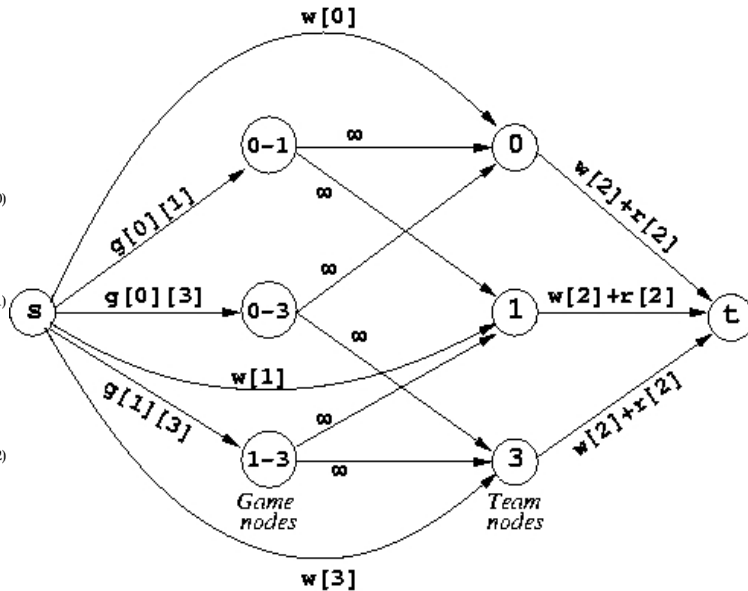
-compute algorithms

-more efficient compute algorithms

-software implementation

-Enterprise strength software implementation

$$\begin{aligned} & \frac{D}{Dt} \overline{w^i w^j} + \overline{w^i w^j} \nabla_x \bar{u}^i + \overline{w^j w^i} \nabla_x \bar{u}^j - \alpha \left( \overline{g^i w^j \frac{T}{T}} + \overline{g^j w^i \frac{T}{T}} \right) \left( \nabla_x \bar{\Phi} + \frac{D\bar{u}_i}{Dt} \right) \\ & + \frac{1}{\rho} \nabla_x [\bar{\rho} u^i w^j w^i + (\bar{g}^i w^j + \bar{g}^j w^i) P^i - \overline{w^i \sigma^i(u)} - \overline{w^j \sigma^j(u)}] \\ & + \frac{1}{\rho} \overline{w^i w^j \nabla_x (\bar{\rho} u^i)} - \overline{P^i (g^i \nabla_x w^j + g^j \nabla_x w^i)} = -\frac{1}{\rho} [\overline{\sigma^i(u) \nabla_x w^j} + \overline{\sigma^j(u) \nabla_x w^i}] = -e_2^j, \quad (30) \\ (1 + e_4) \frac{D}{Dt} \left( \frac{T}{T} \right)^2 - 2f(t) \left( \frac{T}{T} \right)^2 - 2w^a \frac{T}{T} D_x + \frac{1}{(1 + e_4) \bar{\rho} C_p^2} \nabla_x \left[ (1 + e_4) C_p^2 \bar{\rho} w^a \left( \frac{T}{T} \right)^2 \right] + \frac{1 + e_4}{\rho} \left( \frac{T}{T} \right)^2 \nabla_x (\bar{\rho} u^a) \\ & + \frac{2}{\bar{\rho} T C_p} \frac{T}{T} \left[ P^i \nabla_x w^a - \nabla_x (P^i w^a) - \frac{D P^i}{Dt} \right] = \frac{2}{\bar{\rho} T C_p} \frac{T}{T} [\overline{\sigma^i(u) \nabla_x w^a} - \nabla_x F^i] = -e_2, \quad (31) \\ (1 + e_4) \left[ \frac{D}{Dt} \left( \frac{T}{T} \right) + w^a \frac{T}{T} \nabla_x \bar{u}^i - \alpha \left( \frac{T}{T} \right)^2 g^i \left( \nabla_x \bar{\Phi} + \frac{D\bar{u}_i}{Dt} \right) \right] - f(t) w^i \frac{T}{T} - \overline{w^i w^a} D_x \\ & + \frac{1}{\bar{\rho} C_p} \nabla_x \left[ (1 + e_4) C_p \bar{\rho} w^i w^a \frac{T}{T} \right] + \frac{1 + e_4}{\rho} w^i \frac{T}{T} \nabla_x (\bar{\rho} u^a) + \frac{1}{\bar{\rho} T C_p} w^i \left[ P^i \nabla_x w^a - \nabla_x (P^i w^a) - \frac{D P^i}{Dt} \right] \\ & = \frac{1 + e_4}{\rho} \frac{T}{T} \nabla_x \sigma^i(u) + \frac{1}{\bar{\rho} T C_p} \overline{w^i [\sigma^i(u) \nabla_x w^a - \nabla_x F^i]} = -e_2^i, \quad (32) \end{aligned}$$



```

1 /* This line basically imports the "stdio" header file, part of
2 * the standard library. It provides input and output functional
3 * to the program.
4 */
5 #include <stdio.h>
6
7 /*
8 * Function (method) declaration. This outputs "Hello, world" to
9 * standard output when invoked.
10 */
11 void sayHello() {
12     // printf() in C outputs the specified text (with optional
13     // formatting options) when invoked.
14     printf("Hello, world!");
15 }
16
17 /*
18 * This is a "main function". The compiled program will run the
19 * defined here.
20 */
21 void main() {
22     // Invoke the sayHello() function.
23     sayHello();
24 }

```

<http://www.cs.princeton.edu/courses/archive/spr05/cos226/assignments/baseball/>

[http://www.wikiwand.com/en/Programming\\_language](http://www.wikiwand.com/en/Programming_language)

## Mission

## Products

## Awards

## Value

DNaseq

TNseq

### DNaseq

Identical\* results as  
Broad Institute's  
"Best Practice Workflow"  
BWA-GATK HaplotypeCaller

### TNseq

Identical\* results as  
Broad Institute's  
"Somatic Variant Discovery Workflow"  
MuTect and MuTect2

- Identical math, much more efficient computing algorithm and enterprise-strength software engineering
- 10+ X faster whole pipeline FASTQ to VCF; 20X-50X faster GATK/MuTect/MuTect2 portion; in core-hours
- Both products kept up-to-date with Broad Institute's releases

\*1/1000 vcf differences due to GATK down-sampling, thread dependency, rounding differences

DNAseq

TNseq

## GATK is the gold standard

GATK Haplotype Caller is the most accurate DNA analysis tool

- “As of GATK version 3.3, we recommend using HaplotypeCaller in all cases, with no exceptions.”
  - Broad Institute
- “Haplotype Caller is more accurate than the Unified Genotyper”
  - [https://hpc.mssm.edu/files/Carneiro\\_workshop.pdf](https://hpc.mssm.edu/files/Carneiro_workshop.pdf)
- “GATK HaplotypeCaller called a substantial number of indels not called using VarScan-Cons (as well as GATK UnifiedGenotyper)”
  - Warden CD, Adamson AW, Neuhausen SL, Wu X. PeerJ. 2014 Sep 30;2:e600. doi: 10.7717/peerj.600. eCollection 2014.

DNAseq

TNseq

## GATK is the gold standard

GATK Haplotype Caller is the most accurate DNA analysis tool

But GATK Haplotype Caller is too slow

Other speedup efforts:

- (1) massive parallel on cloud (challenges: workflow, data privacy, cost) – may combine with our solution;
- (2) hardware acceleration (challenges: cost, intrusive, inflexibility, scalability);
- (3) corner-cutting (challenge: lose info)

Our approach: stay with the most accurate math, but use much more efficient compute algorithm with enterprise-strength software implementation

## Mission

## Products

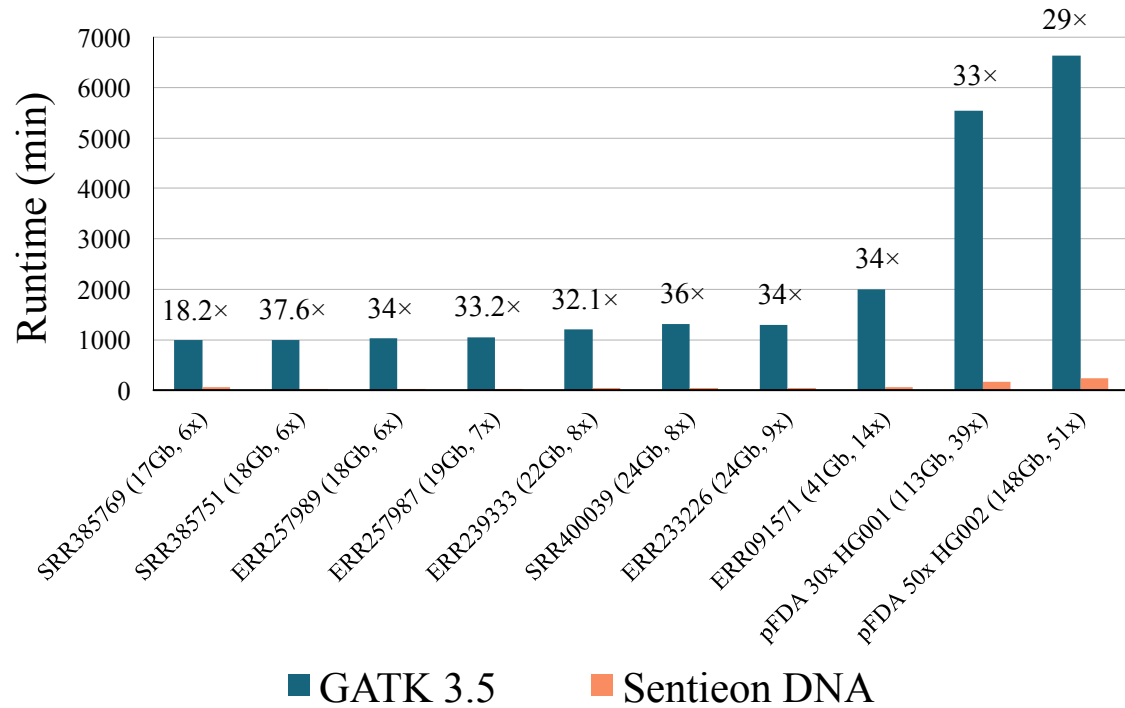
## Awards

## Value

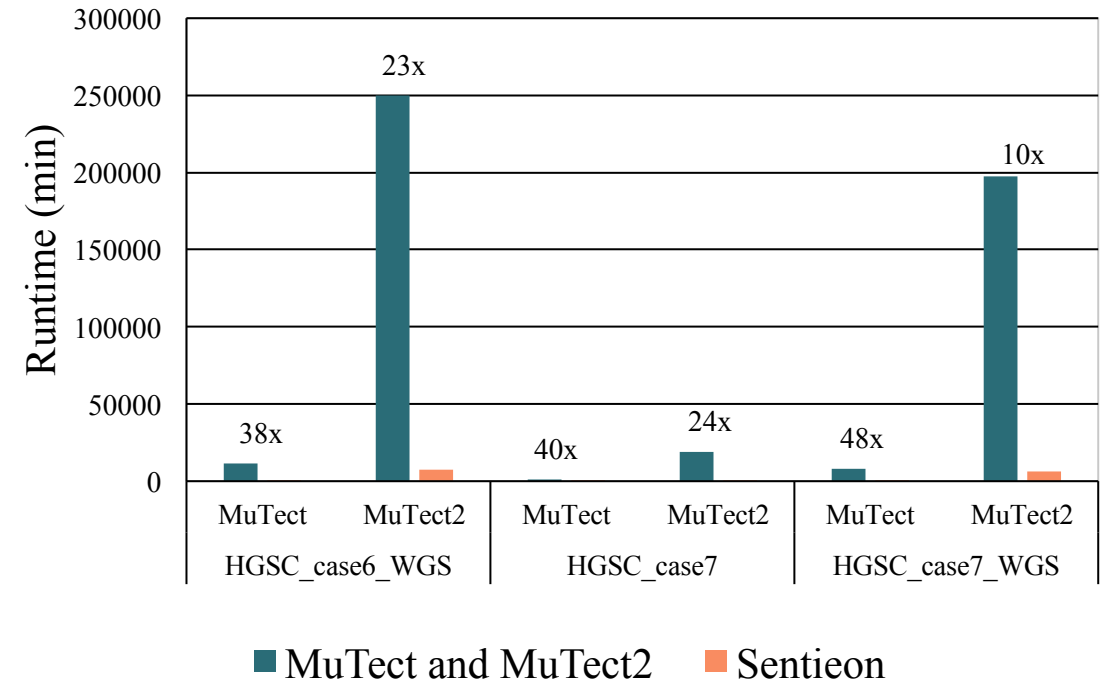
DNaseq

TNseq

### DNaseq



### TNseq



\*Server specs: 32 core 2.4 GHz Intel Xenon server, 64 GB memory

DNAseq

TNseq

## Highlights beyond **speed**:

- 100% **consistency**, no run-to-run differences
  - No down-sampling in high coverage region, no thread dependency
  - ➔ Higher **accuracy** by eliminating software noise
- System **robustness**
  - ➔ Large dataset joint call over 100K samples together (without intermediate merging)

precisionFDA

DREAM

### Consistency Challenge

(2/25-4/25/2016)

- Top Overall Performance
- Highest Reproducibility

### Truth Challenge

(4/26-5/26/2016)

- Highest INDEL Precision
- Highest SNP Recall



<p>HIGHEST Reproducibility</p> <p>AWARDED TO <b>Sentieon team</b></p> <p>Rafael Aldana Hanying Feng Brendan Gallagher Jun Ye</p>	<p>TOP OVERALL Performance</p> <p>AWARDED TO <b>Sentieon team</b></p> <p>Rafael Aldana Hanying Feng Brendan Gallagher Jun Ye</p>
<p>HIGHEST SNP Recall</p> <p>AWARDED TO <b>Sentieon</b></p> <p>Rafael Aldana Hanying Feng Brendan Gallagher Jun Ye</p>	<p>HIGHEST INDEL Precision</p> <p>AWARDED TO <b>Sentieon</b></p> <p>Rafael Aldana Hanying Feng Brendan Gallagher Jun Ye</p>

\*Screenshots from <https://precision.fda.gov/>



## Mission

## Products

## Awards

## Value

precisionFDA

DREAM

ICGC-TCGA DREAM Mutation Calling challenge

---



Annual open contest by ICGC-TCGA for somatic variant calling accuracy.

Challenge-6 due 8/19/2016 (extended from 4/22/2016)

Mission

Products

Awards

Value

precisionFDA

ICGC-TCGA DREAM Mutation Calling challenge

DREAM



Final Leaderboard (8/19/2016)

	SNV	INDEL	SV
Sentieon	98.57%	Sentieon 98.14%	Sentieon 100%
on 4/21	98.17%	on 4/21 97.48%	on 4/21 100%
Bina/Roche	97.57%	Bina/Roche 97.01%	Genowis 99.82%
Genowis	96.92%	OICR-GSI 86.99%	Gridss 99.63%

**-Sentieon leads in all categories-**

## Consistency

precisionFDA Consistency Challenge Reproducibility  
F1-score(%) between runs and between samples

	Garvan vs. Garvan rerun			Garvan vs HLI		
	FP	FN	F1-score	FP	FN	F1-score
Sentieon by UNM*	0	0	100	134633	315831	95.07
Sentieon (dual mapping)	0	0	100	107302	295568	95.53
ISAAC**	0	0	100	112255	266952	95.32
Genalice**	3621	3673	99.92	286306	433504	92.08
Edico Dragen**	3147	3216	99.93	161213	315611	94.87

\* Sentieon standard pipeline results from Jeremy Edwards (University of New Mexico) submission

\*\* Edico, Genalice (MAP 2.2.0) and Isaac (aligner v01.14.07.14 & variant caller v2.0.13) run results from Changhoon Kim's (Macrogen Clinical Laboratory) submission

## Accuracy

precisionFDA Consistency Challenge Accuracy  
F1-score(%) to NIST truth set

	Garvan sample			HLI sample		
	All	SNP only	Indel only	All	SNP only	Indel only
Sentieon by UNM*	99.39	99.86	95.85	98.97	99.73	92.94
Sentieon (dual mapping)	<b>99.47</b>	<b>99.88</b>	<b>96.37</b>	<b>99.06</b>	<b>99.77</b>	<b>93.46</b>
ISAAC**	97.29	98.57	86.66	96.34	97.99	82.10
Genalice**	98.04	99.20	89.10	97.25	98.83	84.66
Edico Dragen**	99.25	99.74	95.49	98.85	99.62	92.74

\* Sentieon standard pipeline results from Jeremy Edwards (University of New Mexico) submission

\*\* Edico, Genalice (MAP 2.2.0) and Isaac (aligner v01.14.07.14 & variant caller v2.0.13) run results from Changhoon Kim's (Macrogen Clinical Laboratory) submission

## Accuracy

precisionFDA Truth Challenge Accuracy  
F1-score(%) to NIST truth set

	HG002 sample			HG001 sample		
	All	SNP only	Indel only	All	SNP only	Indel only
Sentieon	99.8950	99.9548	99.3628	99.8971	99.9605	99.3324
GATK (Sanofi)	99.8905	99.9456	99.4009	99.7053	99.7483	99.3229
Consensus (Bina)	99.8805	99.9382	99.3675	99.8776	99.9325	99.3903
Verily	99.8597	99.9587	98.9802	99.8620	99.9554	99.0328
Edico Dragen (Macrogen)	99.7569	99.8268	99.1359	99.7739	99.8454	99.1397
GATK (Garvan)	99.5679	99.5934	99.3424	99.6208	99.6487	99.3742
Genalice (Macrogen)	98.9492	99.5188	93.9183	98.9624	99.5204	94.0347
Isaac (Macrogen)	98.2618	98.5357	95.8099	98.2800	98.5433	95.9257

### Observations:

- Much higher F1 score (~99.9%) than Consistency Challenge, due to masking out complex regions in truth set
- Sentieon has excellent consistency between HG001 and HG002, due to no run-to-run differences

## Value of Sentieon solutions

- Highest **accuracy**: most rigorous math, no noise in algorithm and software
- **No down-sampling** in high-coverage regions ← Critical for clinical samples
- **No run-to-run difference** ← Critical for medical decision
  - proven in precisionFDA challenges and DREAM challenge
- **Fast turnaround**, ability to scale below 20 minute turnaround → Improved productivity, faster medical decision
- **>10X reduced core-hours** → drastically reduced compute cost
- Enable **large dataset joint call** → Enable genomics big data analysis

*Sentieon tools: fastest, most accurate, zero run-to-run difference, no down-sampling, large cohort joint call, pure software solution, running on any generic computer systems*

## Mission

## Products

## Awards

## Value

DNaseq

TNseq

## Sentieon Software in Use

- Deployed at >100 sites worldwide (academia and industry)
- Accumulated usage at customer sites:  
>100K WGS/WES, >5e15 bases

*Sentieon products are built on the solid foundation of the most rigorous and most extensively validated mathematical models used in Broad Institute's Best Practice Workflow, but with more efficient computing algorithms and enterprise-strength software implementation*

# Thank You

Contact Don Freed ([don.freed@sentieon.com](mailto:don.freed@sentieon.com))  
if you would like to talk to me in person