PRESAGE: PRIVACY-PRESERVING GENETIC TESTING VIA INTEL SGX

Motivation

It is becoming a big challenge to efficiently store and process the huge amount of genomic data for individual biomedical research institutions.

Cloud Computing emerges as an ideal platform for providing elastic computation and storage resources for genomic data analysis





Motivation

 Individual genomic information tends to reveal sensitive personal information. Thus privacy concerns have posed challenges to outsource genomic data in an untrusted cloud environment •*Lin et. al. 2004 Science*: 75 or more SNPs (Single-nucleotide polymorphism) will be sufficient to identify a single person.

•*Gymrek et al. 2013 Science*: surnames can be recovered from personal genomes, linking Utah Residents with Northern and Western European Ancestry (CEU) and public genetic genealogy databases (Ysearch & SMGF).

•*Lipper et.al. 2017 PNAS*: Prediction of human physical traits and demographic information.



Homer et al. 2008 PLOS Genetics: aggregated genome data (i.e., allele frequencies) can also be used for re-identifying an individual in a case group with a certain disease.

Our Solution

We present one of the first implementations of SGX based secure genetic testing framework to facilitate efficiently outsourced storage and computation.

The secure outsource storage is achieved through data sealing scheme within SGX framework, which is immune to replay attack.

We have taken into account the oblivious access protection by using 4KB page-wise data access model.

To improve the performance, we adopt a perfect hashing scheme to achieve O(1) complexity data access within each 4KB page.



Intel SGX Based Secure Genetic Testing Cloud

Our Solution



Workflows of the proposed PRESAGE framework, presented in three consecutive steps:

- 1. Preprocessing
- 2. Encryption and data outsourcing.
- 3. Secure Genetic Query Matching.

Experimental Studies

- Dataset: The dataset is presented in VCF format. And sizes of VCF datasets used in our experiments vary from 10,000 to 200,000 records.
- Experiment Environment: All of the experiments except the iDASH competition results are conducted on a Windows 10 SGX-enabled machine with i7 6820HK CPU and 48 GB memory. Both data owner and CSP were simulated on the aforementioned SGX machine. The iDASH competition results were evaluated on the Linux server with an Xeon Processor E3-1275 v5 and 64 GB memory.

Results

Comparison of querying performance among PRESAGE, HME-based method and plaintext implementation

Table 1. The breakdown run time (in seconds) of the proposed PRESAGE framework

Becord size	Plaintaxt	Encoded data	Socied data	Enclave memory usage			
Record Size	Flaintext		Sealed data	1 query	3 queries		
10,000	0.55 MB	0.09 MB	0.12 MB	3.006 MB	3.006 MB		
50,000	2.59 MB	0.45 MB	0.59 MB	3.010 MB	3.010 MB		
100,000	5.26 MB	0.90 MB	1.15 MB	3.010 MB	3.010 MB		
200,000	10.5 MB	1.75 MB	2.31 MB	3.010 MB	3.010 MB		

Table 2. The data size and enclave memory consumption (in MB) for different datasets.

Record size Attestation	Attoctation	SNDs and ing	Hach concretion	Enclave	Data cooling	Number of queries		
	Allestation	SNPS Coulling	nash generation	creation	Data Sealing	1	3	
10,000	0.121s	0.016s	1.130s	0.169s	0.094s	0.003s	0.003s	
50,000	0.126s	0.080s	6.371s	0.173s	0.517s	0.012s	0.013s	
100,000	0.124s	0.164s	13.473s	0.179s	0.980s	0.023s	0.025s	
200,000	0.120s	0.309s	28.677s	0.171s	2.045s	0.043s	0.048s	

Conclusion

- We proposed a secure outsourcing framework, which can defend malicious attack. To improve the efficiency, an minimal perfect hashing scheme has been incorporated
- Our experiment results demonstrated the efficiency of the proposed PRESAGE framework. For a VCF file with 200K records, the PRESAGE securely processes a query within 0.05 seconds, which includes file loading, unsealing and query matching. Compared with state-of-the-art HME solution, PRESAGE framework shows at least 120x performance gain.

Acknowledgements

- Feng Chen,
- Chenghong Wang,
- Wenrui Dai,
- Xiaoqian Jiang,
- Noman Mohammed,
- Md Momin Al Aziz,
- Md Nazmus Sadat,
- Kristin Lauter,
- Shuang Wang

Secure GWAS via IntelSGX

Motivation / Goal

Enable secure whole genome variant search among multiple individuals from multiple institutions

Institution A has VCF files from x_A individuals labeled case/control, Institution B has VCF files from x_B individuals labeled case/control, ...

Institutes don't want to share data, want to do GWAS on untrusted cloud

Requirements

Secure (Everything kept encrypted outside the SGX Enclave)

Fast

Accurate

Challenges

Data is too large (iDASH dataset ~30GB, real-life much bigger)

SGX Enclave max size 128 MB

Linux allows 4GB paging – paging is extremely slow, typically many orders of magnitude slower than RAM

More data longer transfer, more data slower encryption/decryption

Solution Outline

Keep hash table inside the SGX enclave (server)

Filter and compress VCF files (Client(s))

Construct Enclave (Server)

Perform Remote Attestation (Both Parties Exchange Messages)

Receive Data, Update Hash Table (Server)

Calculate Top-K (K=10 for iDASH) SNPs wrt X² test (Server)

Allele Counting for X² Test

	GG	GT	TT	Total
Cases	<i>r</i> ₀	<i>r</i> ₁	<i>r</i> ₂	R
Controls	s_0	s_1	s ₂	S
Total	<i>n</i> ₀	<i>n</i> ₁	<i>n</i> ₂	N

	Observed allele counts						
	G	т	Total				
Cases	$2r_0 + r_1$	$r_1 + 2r_2$	2R				
Controls	$2s_0 + s_1$	<i>s</i> ₁ +2 <i>s</i> ₂	2 <i>S</i>				
Total	$2n_0 + n_1$	$n_1 + 2n_2$	2N				

Expected allele counts						
G	т					
$2R(2n_0+n_1)/(2N)$	$2R(n_1+2n_2)/(2N)$					
$2S(2n_0+n_1)/(2N)$	$2S(n_1+2n_2)/(2N)$					

Chi-square test for independence of rows and columns (null hypothesis):

$$\sum \frac{(Obs - Exp)^2}{Exp} \sim \chi^2$$
 with 1 df

VCF (Variant Call Format)

##real	id in 10	00genome	project:	HG03518	3			
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	TYPE	
1	13380	rs571093	408	С	G	100	PASS	heterozygous
1	15211	rs786018	09	Т	G	100	PASS	heterozygous
1	15820	rs269131	5	G	Т	100	PASS	heterozygous
1	18849	rs533090	414	С	G	100	PASS	heterozygous
1	30923	rs806731		G	Т	100	PASS	heterozygous
1	49298	rs200943	160	Т	С	100	PASS	heterozygous
1	52238	rs269127	7	Т	G	100	PASS	heterozygous
1	55164	rs309127	4	С	A	100	PASS	heterozygous
1	62777	rs528401	309	A	Т	100	PASS	heterozygous
1	69897	rs200676	709	Т	С	100	PASS	heterozygous
1	82343	rs563238	524	Т	С	100	PASS	heterozygous
1	83084	rs181193	408	Т	A	100	PASS	heterozygous
1	86331	rs115209	712	A	G	100	PASS	heterozygous

Filtering & Variable Length Encoding

[203760]	2182	\rightarrow	100010000110	12
[203761]	11414	\rightarrow	10110010010110	14
[203762]	2463	\rightarrow	100110011111	12
[203763]	2083	\rightarrow	100000100011	12
[203764]	1169	\rightarrow	10010010001	11
[203765]	20377	\rightarrow	100111110011001	15
[203766]	7717	\rightarrow	1111000100101	13
[203767]	15460	\rightarrow	11110001100100	14
[203768]	5258	\rightarrow	1010010001010	13
[203769]	8842	\rightarrow	10001010001010	14
[203770]	2846	\rightarrow	101100011110	12

Filtering VCF and Compressed Representation

- Only SNP "ID" and "TYPE" columns essential
- "QUAL" and "FILTER" can be removed during preprocessing
- "CHROM", "POS", "REF", "ALT" can all be found via "ID" from dbSNP
- Trim the "rs" in front of "ID", represent as integer
- Sort by "TYPE", so that we don't have to keep heterozygous/homozygous
- Keep a single integer to determine when "TYPE" changes

Variable Length Encoding

- ► Sort "ID"s, grouped by "TYPE"
- Keep only differences using the minimum number of bits needed
- Keep another small stream for bit-lengths, encoded by a Huffman Tree
- Example VCF Filtering/Compression: Actual VCF Size: 15,428,390 Bytes
- Heterozygous-Stream: 944,166 bits Homozygous-Stream: 428,921 bits Total Main Stream Size: 1,373,087 bits
- 5-bits/len Auxiliary Stream Size: 1,631,105 bits Huffman Encoded Auxiliary Stream Size: 1,070,458 bits
- Main Stream + Huffman Auxiliary Stream: 305,444 Bytes

Preliminary Results

1000 case / 1000 control. ~300K-350K SNPs per VCF, ~5.5M unique SNPs

SGX Enclave creation: 0.193381 seconds

Remote Attestation: 0.002464 seconds

Main Application: 49.990706 seconds

SGX Enclave destruction: 0.034888 seconds

Acknowledgements

- Can Kockan (Indiana University)
- Natnatee Dokmai (Indiana University)
- Oguzhan Kulekci (Istanbul Technical University)
- Steve Myers (Indiana University)
- The Cancer Genome Collaboratory
- Indiana University Precision Health Initiative

Precision genomics under scalability, privacy and security constraints

S. Cenk Sahinalp Indiana University Bloomington Vancouver Prostate Centre

Genome storage and communication: the need

- Research: massive genome projects (e.g. PCAWG) require to exchange 10000s of genomes.
- Need to cover >200 cancer types, many subtypes.
- Within PCAWG TCGA to sequence 11K patients covering 33 cancer types. ICGC to cover 50 cancer types >15PB data at The Cancer Genome Collaboratory
- Clinic: \$100s/genome (e.g. NovaSeq) enable sequencing to be a standard tool for pathology
- PCAWG estimates that 250M+ individuals will be sequenced by 2030

Current needs

- Typical technology: Illumina NovaSeq 100-400 bp reads, 250-500GB uncompressed data for high coverage human genome, high redundancy
- 40% of the human genome is repetitive (mobile genetic elements, centromeric DNA, segmental duplications, etc.)
- Upload/download: 55 hrs on a 10Mbit consumer line; 5.5 hrs on a 100Mbit high speed connection
- Technologies under development: expensive, longer reads, higher error (PacBio, Nanopore) – lower redundancy and higher error rates limit compression

File formats

- Raw read data: FASTQ/FASTA dominant fields: (read name, sequence, quality score)
- Mapped read data: SAM/BAM reads reordered based on mapping locus on a reference genome (not suitable for metagenomics, organisms with no reference)
- Key decisions to be made:
- Each format can be compressed through specialized methods should there be a standardized format for compressed genomes?
- Better file formats based on mapping loci on sequence graphs representing common variants in a pan-genomic reference?

Genome Compression: Towards an International Standard

- Collaboration with MPEG to evaluate the current state of HTS data compression towards an International Standard
- Standard Benchmark DataSet: 2+ TB sequence data:
 - 7 FASTQ samples and 8 SAM samples, covering 6 species, 6 technologies, various use-cases (high and low coverage data, cancer cell lines, WGS, RNA-Seq, metagenomics etc.)
- 15 FASTQ tools and 10 SAM tools evaluated
- Available at https://github.com/sfu-compbio/compression-benchmark

Comparison of high-throughput sequencing data compression tools [Numanagić et al., Nat. Meth., Dec 2016]

FASTA/Q compression

General purpose compressors used in genomics

- LZ77 tools (gzip, pigz)
- BWT tools (bzip2, pbzip2)
- LZMA (7z)
- Context mixing (zpaq, lpaq)
- NCBI Toolkit (used at SRA for storing samples)

 General compressors do not take into account redundancies specific to FASTQ format (FASTQ files are treated as ordinary text files)

Compressor (on 53.8 GB human g. 6.5x coverage)	Size (total)	Size (field by field)	Size (sequence)
pigz	18.5 GB	16.1 GB	5.9 GB
pbzip2	14.8 GB	14.1 GB	5.4 GB
NCBI SRA	~ 14.2 G	B	

Specialized FASTA/Q compressors

• Goals:

- Read name tokenization
- Separate sequence and quality score modeling
- Examples:
 - DSRC and DSRC2 [1] (Huffman coding)
 - fastqz, fqzcomp [2] and Slimfastq [3] (context mixing with arithmetic coding)
 - FQC [4] and LFQC [5] (LZMA, paq and ppmd as compression engine)

Compressor (on 53.8 GB)	Size (total)	Size (sequence)		
DSRC2	13.2 GB	5.2 GB		
Slimfastq	11.0 GB	4.4 GB		
FQC	11.4 GB	N/A		

Roguski S, Deorowicz S. DSRC 2--Industry-oriented compression of FASTQ files. Bioinformatics, 2014
 Bonfield JK, Mahoney MV. Compression of FASTQ and SAM Format Sequencing Data. PLoS ONE, 2013
 Ezra J. <u>https://github.com/Infinidat/slimfastq</u>

[4] Dutta A, Haque MM, Bose T, Reddy CV, Mande SS. FQC: A novel approach for efficient compression, archival, and dissemination of FASTQ datasets. J Bioinform Comput Biol., 2015
[5] Nicolae M, Pathak S, Rajasekaran S. LFQC: a lossless compression algorithm for FASTQ files. Bioinformatics, 2015

FASTA/Q compressors based on read reordering

• Goals:

- Reorder reads to improve locality of reference

Compressor (on 53.8 GB)	Size (total)	Size (sequence)		
SCALCE	10.8 GB	3.0 GB		
ORCOM	N/A	1.7 GB		
Mince	N/A	6.0 GB		
LW-FQZip	N/A			

- Examples:
 - SCALCE [1] (uses locally consistent parsing for read reordering/clustering)
 - ORCOM [2] (uses lexicographically smallest kmers for clustering)
 - Mince [3] (similar to ORCOM)
 - LW-FQZip [4] (uses implicit mapping to a reference)

[1] Hach F, Numanagić I, Alkan C, Sahinalp SC. SCALCE: boosting sequence compression algorithms using locally consistent encoding. Bioinformatics, 2012

[2] Grabowski S, Deorowicz S, Roguski L. Disk-based compression of data from genome sequencing. Bioinformatics, 2014

[3] Patro R, Kingsford C. Data-dependent Bucketing Improves Reference-free Compression of Sequencing Reads. Bioinformatics, 2015

[4] Zhang Y, Li L, Yang Y, Yang X, He S, Zhu Z. Light-weight reference-based compression of FASTQ data. BMC Bioinformatics, 2015

FASTA/Q compressors based on read assembly

• Goals:

- Assemble the underlying genome and map reads to the assembly

- Examples:
 - Quip [1] (Bloom filters, assembles clusters of 1 million reads)
 - Leon [2] (probabilistic de Bruijn graph)
 - k-Path [3] (probabilistic de Bruijn graph)

Compressor (on 53.8 GB)	Size (total)	Size (sequence)			
Quip	11.3 GB	4.5 GB			
Leon	13.6 GB	4.7 GB			
k-Path	N/A	2.0 GB			

[1] Jones DC, Ruzzo WL, Peng X, Katze MG. Compression of next-generation sequencing reads aided by highly efficient de novo assembly. Nucleic Acids Res. 2012

[2] Benoit G, Lemaitre C, Lavenier D, Drezen E, Dayris T, Uricaru R, Rizk G.. Reference-free compression of high throughput sequencing data with a probabilistic de Bruijn graph. BMC Bioinformatics, 2105.
 [3] Kingsford C, Patro K. Reference-based compression of short-read sequences using path encoding. Bioinformatics, 2015

Compression results on raw (FASTA/Q) read data

Sample	SRI	2554369	SRI	R327342	MH000	1.081026	SRR1284073		SRR870667		SRR870667 ERR174310		ERR174324														
Organism	P.ae	eruginos a	<i>S</i> .	cerevisiae	H.sa	piens Gut		E.coli	T.cacao		H.	sapiens	H.,	sapiens													
Technology	Illumir	na GAIIx	Illum	ina GAII	Illu	Illumina GA		Illumina GA Pach		PacBio		PacBio		PacBio Illumina GAIIx		ia GA PacBio Illumina GAIIx Hi		PacBio Illumina GAIIx		Illumina GAIIx		Illumina GAIIx HiSeq		Illumina GAIIx			HiSeq
Coverage		105x	1	Unknown		Unknown		5x		65x		25x		335x													
0-1-11	550		3,881		1,880		1,309		22,944		53,869		2,717,029														
Original	165		947		512		649		7,463		20,966		1,059,387														
plag	158	1.00	1,020	1.00	501	1.00	547	1.00	6,943	1.00	18,597	1.00	305,690	1.00													
pigz	48	1.00	277	1.00	149	1.00	188	1.00	2,108	1.00	5,982	1.00	104,927	1.00													
phrip?	125	1.19	831	1.41	390	1.33	463	0.74	5,577	1.07	14,887	0.80	242,834	0.21													
pozipz	44	6.12	251	6.80	139	6.10	176	6.99	1,879	3.59	5,473	3.08	95,969	1.23													
DSPC2	105	0.21	668	0.26	312	0.25	N	/ •	4,761	0.23	13,214	0.20	NI/A														
DSILOZ	41	2.15	257	3.22	128	2.06	14/	/ A	1,865	1.45	5,239	1.25	N/A														
Fazomp	89	0.35	559	0.37	280	0.41	N	/ ^	4,028	0.34	11,320	0.31	N/A														
rqzcomp	37	N/A	203	7.39	120	N/A	14,	A	1,556	N/A	4,623	3.38	N/A														
Fastaz	N	14	N	1	N	/ •	N	/ •	N/	٨	10,955	3.45	NI/A														
rastqz		/A	14,	A	14	/ A	14,	A	14/1	^	N/A	N/A	N/A														
Slimfasta	94	0.54	507	0.48	266	0.54	N	/Δ	4,280	0.52	11,045	0.47	178,092	0.49													
Similasty	30	11.62	149	9.93	104	10.94		A	1,416	5.82	4,426	4.89	77,629	5.94													
FOC	76	1.04	494	1.23	268	1.51	413	0.98	3,912	1.20	11,409	1.22	N/A														
140	N/A	12.16	N/A	13.42	N/A	18.66	N/A 12.12		N/A	6.34	N/A	5.87	N/A														
LEOC	69	9.24	490	8.67	266	10.44	407	18.03	2,412	8.53	N	4	N/A														
Lido	17	159.86	129	146.15	103	162.94	156	386.25	N/A	N/A		•	N/A														
SCALCE	76	0.38	487	0.29	297	0.40	421	0.67	3,699	0.35	10,827	0.30	161,067	0.57													
Jonizoz	17	4.59	68	3.87	71	5.83	161	9.78	998	2.88	3,017	2.36	28,452	1.94													
LW-FOZin	117	1.10	790	0.60	N	/ 4	N	N/A 5,038 2.16 N/A		N/A																	
Litt i drub	45	5.60	320	5.25		/			1,735	2.56		•		_													
Quin	89	0.39	537	0.48	272	0.49	420	0.36	3,914	0.50	11,312	0.46	184,051	0.38													
quip	37	7.58	181	8.51	114	11.00	159	10.59	1,462	5.74	4,556	5.39	79,771	4.64													
Leon	87	3.61	544	2.66	291	3.66	479	2.81	4,518	3.93	13,623	3.55	220,397	1.13													
Deon	19	16.95	89	17.02	87	14.22	170	34.31	1,360	10.49	4,739	9.90	83,539	4.66													
KIC	95	5.75	613	7.40	307	5.32	451	9.40	4,498	6.74	13,006	6.19	N/A														
mo	32	6.89	188	7.88	122	7.38	168	9.37	1,594	3.60	4,915	3.30	, M/M														
		0.25		0.22		0.43	N			0.49		0.33		0.12													
Orcom	11	0.77	36	0.43	51	0.90	N,	/A	825	0.72	1,798	0.43	6,921	0.23													
DDDD		4.09		3.14		2.46				3.97		4.39															
BEETL	23	37.52	117	32.22	114	29.83	N,	A	1,200	20.68	3,912	21.76	N/A														
h D-th		1.01		0.81		6.47				1.35	21	1.62	NT / A														
K-Path	14	15.25	45	9.49	62	71.69	IN,	A	660	9.05	2,088	8.55	N/A														

SAM/BAM compression

General purpose compressors for SAM files

- LZ77 tools (gzip, pigz)
- BWT tools (bzip2, pbzip2)
- Current standard: LZ77-based BAM (Samtools [1], Sambamba [2], Picard [3])
- None of those methods treat differently separate SAM columns.
- Clearly, simple stream separation without any additional post-processing increases significantly the overall compression rate

 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, 1000 Genome Project Data Processing Subgroup. The Sequence Alignment/Map format and SAMtools. Bioinformatics, 2009
 Tarasov A, Vilella AJ, Cuppen E, Nijman IJ, Prins P. Sambamba: fast processing of NGS alignment formats. Bioinformatics, 2015
 Broad Institute. http://broadinstitute.github.io/picard/

Compressor (on human cancer g. sample 427 GB)	Size	Size (separate streams)		
pigz	119 GB	103 GB		
pbzip2	100 GB	94 GB		
Samtools	131 GB	102 GB		

Specialized SAM tools

- Separate fields into different compression streams
- Use reference to store sequence information, if possible

- Primarily reference based:
 - CRAM format (Scramble [1], Cramtools [2])
- Assembly and reference based:
 - Quip [3]
- Statistical modeling and arithmetic encoding:
 - sam_comp [4]

Compressor (on human cancer sample; 427 GB)	Size
Cramtools	95 GB
Scramble	82 GB
Scramble (without reference)	86 GB
Quip (without reference)	98 GB
<pre>sam_comp * does not support all SAM fields</pre>	42 GB*

 Bonfield JK.The Scramble conversion tool. Bioinformatics, 2014
 Hsi-Yang Fritz M, Leinonen R, Cochrane G, Birney E. Efficient storage of high throughput DNA sequencing data using reference-based compression. Genome Res., 2011

 [3] Jones DC, Ruzzo WL, Peng X, Katze MG. Compression of next-generation sequencing reads aided by highly efficient de novo assembly. Nucleic Acids Res. 2012
 [4] Bonfield JK, Mahoney MV. Compression of FASTQ and SAM Format Sequencing Data. PLoS ONE, 2013

Local assembly based SAM tools

- Avoid redundant storing of SNPs and other small SVs
- Find the underlying genome via local assembly, and encode SNPs and small SVs only once
- Examples:
 - DeeZ [1]
 - CBC [2]

Compressor (on human cancer sample; 427 GB)	Size
DeeZ	78 GB

Sequence only without assembly	Sequence only with assembly			
4,169 MB	4,120 MB			

 Hach F, Numanagić I, Sahinalp SC. DeeZ: reference-based compression by local assembly. Nat. Methods, 2014
 Ochoa I, Hernaez M, Weissman T. Aligned genomic data compression via improved modeling. Journal of bioinformatics and computational biology, 2014 DeeZ: DeeNA Zeep

Motivation

- BAM (the most common format for storage and communication) misses some opportunities in SAM format, particularly common SNV loci in reads
- Alternative SAM/BAM compression tools, based on arithmetic coding (AC) and other data modeling methods, like Quip and Samcomp, provide superior compression, but not random-access capability
- DeeZ locally assembles reads and represents each SNV once, on the contig.

DeeZ: Quality scores

- Quality scores account for majority of the space in almost any format
 - minor improvement in quality score compression is more beneficial than improvement in other areas

DeeZ on human cancer sample 427 GB	Size	Gain
Sequence only without local assembly (5% of compressed file)	4,169 MB	
Sequence only with local assembly (5% of compressed file)	4,120 MB	49 MB
Quality scores only with order-1 AC model (42% of compressed file)	33,516 MB	
Quality scores only with sam_comp model (41% of compressed file)	31,010 MB	2,506 MB

Compression results on mapped (SAM/BAM) read data

Sample	1	DH10B	9827.2.49 sample-2-1		K562.LID8465		HCC1954		NA12878.S1			
Organism		E.coli	H.sapiens		H.sapiens		H.sapiens		H.sapiens		H.sapiens	
Technology	MiSeq		HiSeq		IonTorrent		RNASeq		Cancer Cell		HiSeq	
Coverage		490x		2x		0.7x		7x		35x		60x
Original	5,579		21,059		5,924		75,915		427,028		589,083	
_:	1,336	0.77	6,021	1.55	1,378	1.48	12,785	1.06	119,839	1.40	113,462	0.13
pigz		0.63		0.82		0.49		0.70		0.91		0.60
phrin?	1,074	1.65	5,243	1.93	1,127	4.04	10,251	3.57	100,280	1.62	89,598	0.46
pozipz		3.16		3.39		3.72		2.46		3.23		0.59
Samtoola	1,407	1.00	6,499	1.00	1,469	1.00	13,757	1.00	131,566	1.00	121,710	1.00
Samtools		1.00		1.00		1.00		1.00		1.00		1.00
Biggard	1,425	1.42	6,517	1.04	1,474	1.82	13,818	1.48	132,861	1.18	NI / A	
Ficard		2.76		1.52		2.10		2.44		1.91	N/A	
Sambamba	1,407	1.05	6,499	0.93	1,469	1.12	13,757	1.05	131,566	1.39	121,710	0.13
Sambamba		1.08		1.13		0.97		0.97		1.12		0.53
Grantaala	1,066	0.93	3,778	1.42	1,170	2.12	10,344	1.70	95,442	1.28	NI / A	
Cramtools		1.71		1.67		4.93		2.00		1.50	N/A	
Saramble	863	0.23	3,297	0.29	1,030	0.62	9,261	0.38	82,041	0.27	66,632	0.10
Scramble		0.76		0.66		1.58		0.67		0.71		0.50
Samemble without reference	899	0.29	4,236	1.18	1,113	0.45	9,839	0.43	86,914	0.37	72,407	0.10
Scramble without reference		0.74		0.63		1.06		0.78		0.79		0.47
Scramble with brin?	851	0.76	3,262	0.62	998	1.50	8,611	1.27	80,094	0.60	N/A	
Scramble with bzipz		0.89		0.66		1.72		0.81		0.82	N/A	
DeeZ	823	0.56	3,221	0.78	1,028	1.81	8,120	0.92	78,473	0.91	62,966	0.26
Deez		3.90		2.46		5.51		3.35		2.94		1.00
DeeZ with brin?	730	0.91	2,734	1.23	918	3.49	7,266	2.01	74,509	1.66	$53,\!497$	0.41
Deez with bzipz		10.11		5.60		9.86		7.91		6.39		1.90
TSC	1,105	2.21	7,939	0.80	1,193	2.55	20,864	3.17	164,627	0.50	NI / A	
150		9.05		2.24		6.75		6.27		2.65	N/A	
Quip	1,103	0.67	4,419	0.94	1,230	0.96	11,186	1.19	98,303	0.83	97,165	0.44
		10.69		7.81		3.37		8.27		9.05		2.18
Quip with reference	803	0.67	NI / /	\	NI /	۸	8,743	1.17	N/A		64,493	0.43
Quip with reference		10.06	11/1	1	19/.	~		8.20	N/A			2.20
sam comp	700	0.68	2,649	0.76	891	1.20	7,023	0.71	42,522	0.62	53,263	0.37
sam_comp		3.36		2.95		6.54		3.56		3.25		2.00

Optimal Compressed Representation of High Throughput Sequence Data via Light Assembly

Cenk Sahinalp

Based on joint work with Kaiyuan Zhu, Tony Ginart, Joseph Hui, Ibrahim Numanagić, Thomas Courtade, David Tse

Current FASTQ Compression Schemes

- General purpose compressors (FASTQ files are treated as ordinary text files)
 - gzip (parallel gzip--pigz), bzip2 (parallel bzip2--pbzip2)

Alignment reference g • use de-no	Compressor (on 53.8 GB human with 6.5x coverage)	Size (total)	Size (field by field)	Size (sequence)	e underlying raph
Quip,use an us	pigz	18.5 GB	16.1 GB	5.9 GB (2.251)	ng
• LW-FC	pbzip2	14.8 GB	14.1 GB	5.4 GB (2.060)	reordering of
the reads	SRA	~ 14.2 G	В		ly boost the

compression rates while avoiding information loss.

• SCALCE, Orcom, Mince

Assembltrie: Our New Compressed Representation

- Combine the advantages of reordering and alignment based compressors.
- The input reads are organized into a *forest* of compact trie-like data structures called **read forest**.
 - Each node v represents a read (a string
 - Each directed edge (u, v) represents covered by its parent, v
 - May contain a single cycle acting as the An example trie-like structure

Combinatorial Optimization Formulation

• Among all possible read forests, our objective is to find the one contains minimum number of symbols, i.e. τ : a trie in the forest T,

$$T^* = \arg \min_{T} \sum_{\tau \in T} \sum_{v \in V_{\tau}} w[v, \pi(v)] \checkmark$$

- The greedy algorithm to build the desired read for $e_{ant}^{shortest}$ suffix of v that
 - Pick for each read u an already processed read v with miffimum w[u, v], set its parent $\pi(u)$ to v $\pi(v)$: the parent node of
 - Identify eather already processed ready a within $m_{0.9}$ with only u (can assume $\pi(u) = NIL$)

<u>Theorem</u>: The greedy algorithm computes the optimal T^* with minimum overlap K.

with the corresponding

w[v, u]: the length of the

vertex set V_{τ} ;

Information Theoretic Upper Bound for HTS Data Compression

- Assembltrie achieves combinatorial optimality
 - For any finite collection \mathcal{R} of reads to be compressed with any explicit or implicit (overlap graph) assembly based compressor, it produces the smallest number of symbols to be encoded for reads.
- Is it possible to obtain better compression performance by a fundamentally different data structure (i.e. representation of reads)?
 - NO. The minimum number of bits needed by any algorithm to describe the reads *R* is given by *H*(*R*).

 $H(\mathcal{R}) \approx NL \log(3)h_2(p) + |G| \cdot H (\text{Poisson}(N/|G|)) + LZ(G)$

- \neq Optimal compression in practice
 - The proof does r ot consider read errors
 - Need to Sequencing errors in Read sampling process etc. Reference

 $h_2(p) = -p \log p - (1-p) \log 1 - p$

genome

Compression Performance (8 Threads, in bit per base)

MPEG HTS FASTO Datasetevisiae

Sample	Read L. / Cov.	Assembltrie	Orcom	Mince	K-Path	SCALCE
P.aeruginosa	100 / 25	0.345	0.518	0.484	0.673	0.821
S.cerevisiae	63 / 80	0.271	0.304	0.312	0.384	0.578
H.sapiens gut	44 / NA	0.757	0.804	0.786	2.545	1.104
T.cacao	108 / 20	1.733	0.884	0.735	0.707	1.070
Sim. T.cacao	108 / 19	0.479	0.667	N/A	N/A	N/A
H.sapiens 1	101 / 7	0.570	0.686	0.746	0.797	1.151
H.sapiens 2	101 / 20	0.322	0.364	N/A	N/A	N/A
			overage			

Ref: Numanagić, I., Bonfield, J. K., Hach, F., Voges, J., Ostermann, J., Alberti, C., ... & Sahinalp, S. C. (2016). Comparison of high-throughput sequencing data compression tools. *Nature methods*.

Running Time (8 Threads, in seconds)

Default running time to generate the compression rates in MPEG

Assembltrie's Performance vs Information theoretic upper bound on compression

Acknowledgements

- •
- National Science Foundation (NSF) CCF-1619081, CCF-1528132 and CCF-0939370 (Center for Science of Information)
- National Institutes of Health (NIH) GM108348
- The Cancer Genome Collaboratory
- Indiana University Precision Health Initiative

