

長短兼備：

高通量DNA定序技術於生物醫學研究之應用 Short and Long Read Sequencing Technologies and their Applications in Biomedical Research

黃憲達 (Hsien-Da Huang)

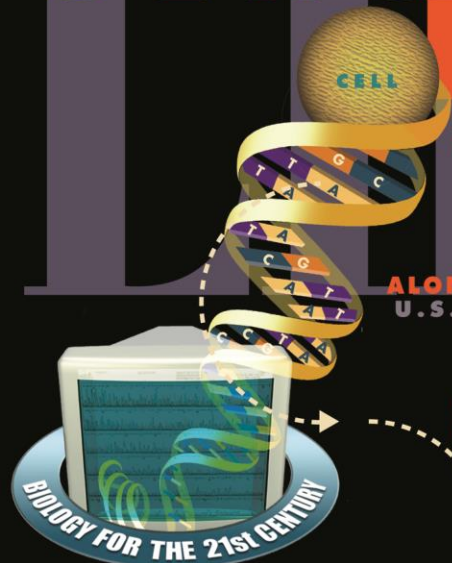
國立交通大學 生物科技學系 講座教授
國立交通大學 生物科技學院 副院長



GENOMES to LIFE

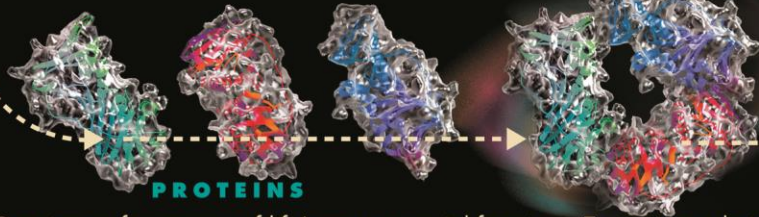
BIOLOGICAL SOLUTIONS FOR ENERGY CHALLENGES

INNOVATIVE APPROACHES ALONG UNCONVENTIONAL PATHS
U.S. DEPARTMENT OF ENERGY



DNA SEQUENCE DATA FROM GENOME PROJECTS

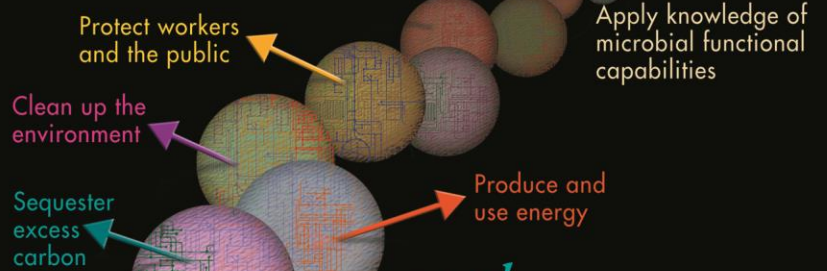
Genes and other DNA sequences contain instructions on how and when to build proteins



PROTEINS

Proteins perform many of life's most essential functions. To carry out their specific roles, they often work together in the cell as protein machines.

goal
IDENTIFY PROTEIN MACHINES



Clean up the environment

Sequester excess carbon

Protect workers and the public

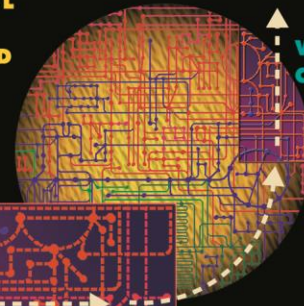
Produce and use energy

Apply knowledge of microbial functional capabilities

goal
EXPLORE FUNCTION IN MICROBIAL COMMUNITIES

COMMUNITY OF CELLS

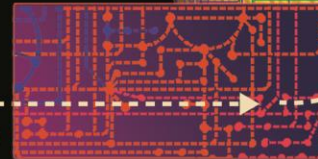
goal
DEVELOP COMPUTATIONAL CAPABILITIES TO UNDERSTAND COMPLEX BIOLOGICAL SYSTEMS



WORKING CELL

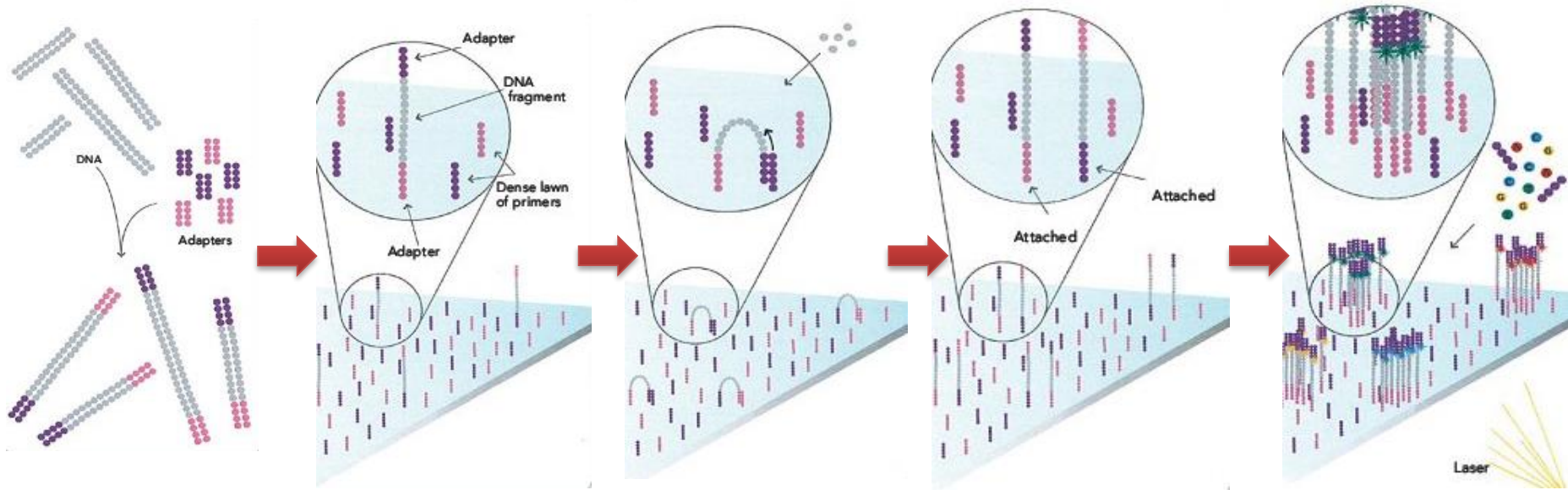
Many protein machines interact through complex, interconnected pathways. Analyzing these dynamic processes will lead to models of life processes.

goal
CHARACTERIZE GENE REGULATORY NETWORKS



URL DOEGenomesToLife.org

Next-generation sequencing



Flow cell



Flow cell

CGATGCGAATT



GCGCAGCGCGGCCGCGCAGCAGCCTCCGCCCCCGCACGGTGT
GTCCCGAGCTAGCCCCGGCGGGCCGCCGCCGCCAGACCGGA
AGTCCCCGCCTCGCCGCCAACGCCACAACCACCGCGCACGG
AGAGCCGGAGCGAGCTCTTCGGGGAGCAGCGATGCGACCCT
CTGCTGGCTGCGCTCTGCCCGGCGAGTCGGGGCTCTGGAGGA
ACGGCTCGTGCGTCCGAGCCTGTGGGGCCGACAGCTATGAGA
GAAGTGCGAAGGGCCTTGCCGCAAAGTGTGTAACGGAATAG
ATAAATGCTACGAATATTAACAACACTTCAAAAACCTGCACCTC
TGGCATTTAGGGGTGACTCCTTCACACATACTCCTCCTCTG
CGTAAAGGAAATCACAGGGTTTTTGGCTGATTCAGGCTTGGC
GAGAACCTAGAAATCATACGCGGCAGGACCAAGCAACATGG
ACATAACATCCTTGGGATTACGCTCCCTCAAGGAGATAAGT
AAATTTGTGCTATGCAAATACAATAAACTGGAAAAAACTGT
ATAAGGAAAGAGAGCTGAAAAGAGCTGGAAAGCCGAGAGCCGA

Sequencing Length, Throughput & Quality

Sanger sequencing



Next-generation sequencing (NGS): illumina platform, Thermo Fisher



3rd or 4th generation sequencing: Pac Bio Science, Oxford Nanopore



豪傑使長槍、君王用短劍

「十八般武器」是傳統兵器總稱。各有形製、功能不同，如「刀」下就有九環刀、雁翎刀、青龍偃月刀等。



武器使用的訣竅，依不同功能而設計不同，如「槍」用於遠刺；「劍」以刺、割為主。

楊鐵心-耍長槍



七十二路「楊家槍法」，出槍長，且虛實，有奇正，進其銳，退其速，其勢險，其節短，穩如山，動如雷。

「十八般武器」分成「九長、九短」，「九長」是：刀、槍、棍、鉞、叉、鑊、鉤、槊、戟。「九短」是：刀、劍、鞭、槁、拐、斧、棒、鎚、杵。



曹操-青釭劍



劉備-雙股劍



孫權-青冥劍

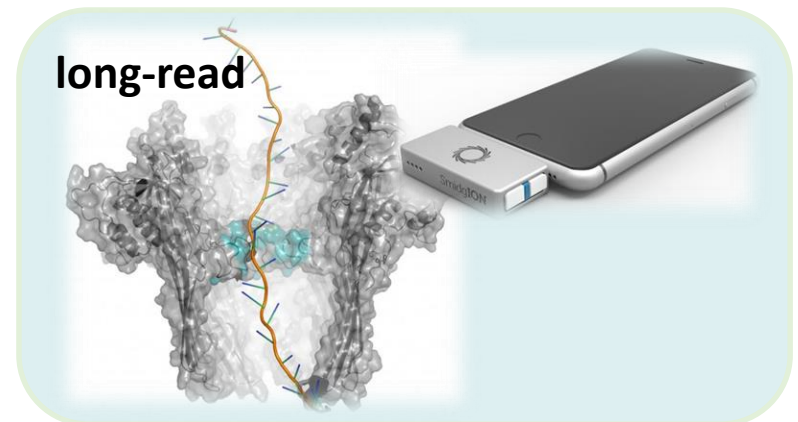
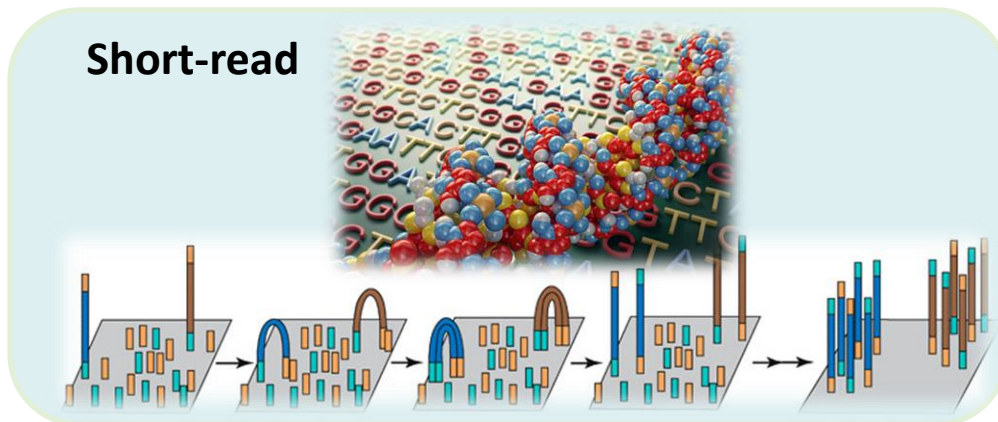


Short-read vs. long-read sequencing

Which one is better?



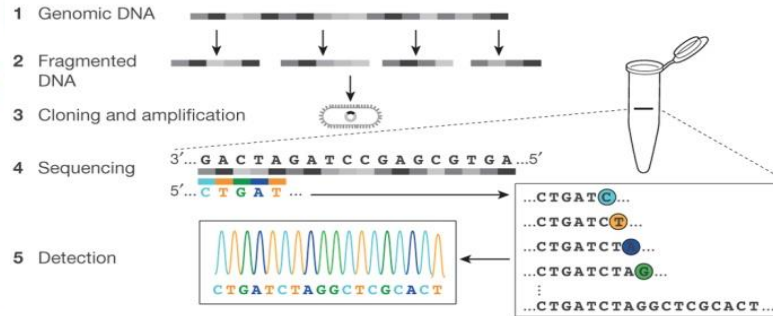
- The **most frequently asked questions** in sequencing.
- The sequencing read length depends on the instrument and chemistry used.
- The range of the read length of a **short-read** sequencing instrument is between **100** and **600 bps**, while that of a **long-read** sequencing instrument varies between **10** to **15 kbps**.
- The choice you make depends on the **goal of your experiment**; one isn't considered universally superior to the other.



DNA sequencing technologies and applications

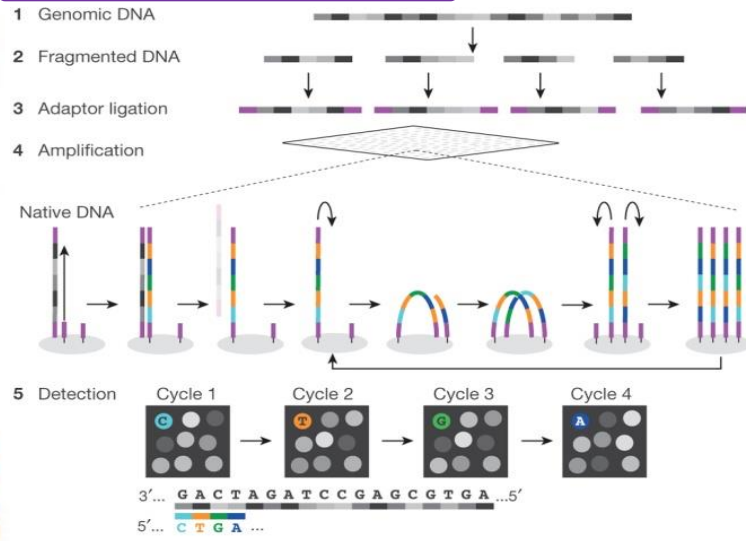
1

First generation sequencing (Sanger)



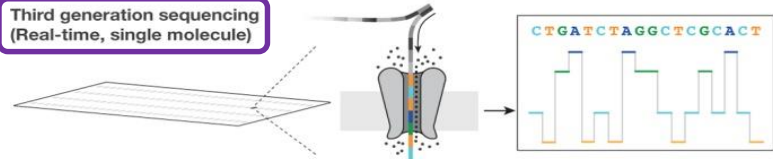
2

Second generation sequencing (massively parallel)

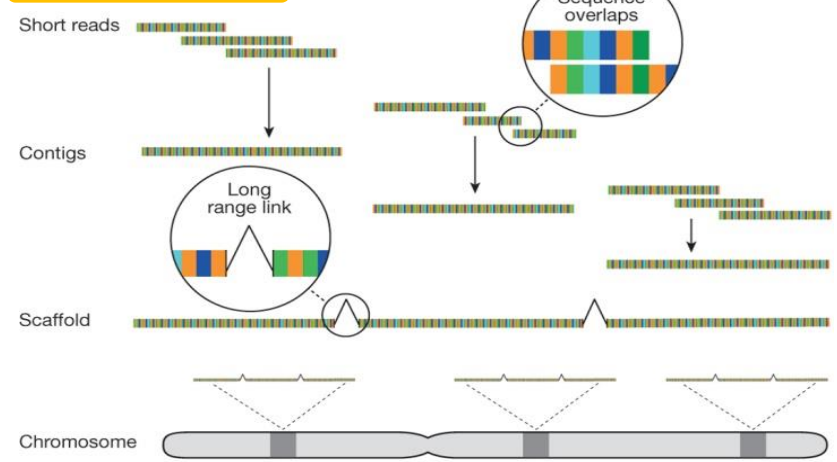


3

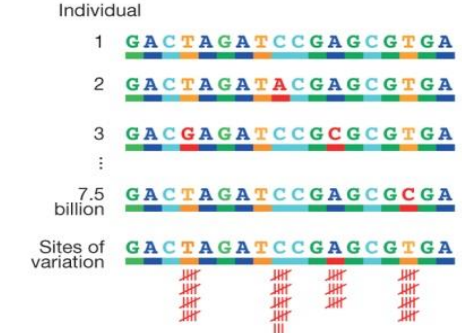
Third generation sequencing (Real-time, single molecule)



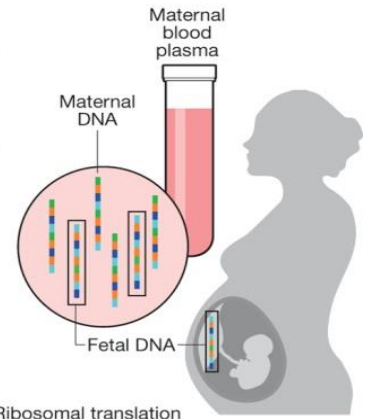
De novo genome assembly



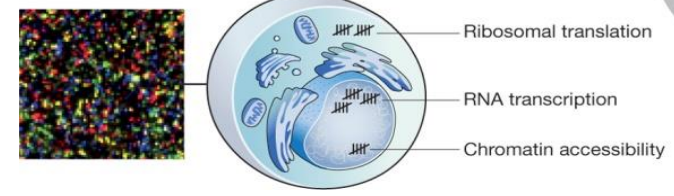
Genome resequencing



Clinical applications (NIPT)



Sequencers as counting devices



Applications of Next-generation Sequencing

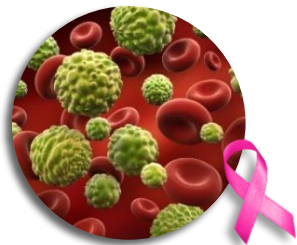
Human Genetics Research



Forensic Genomics



Complex Disease Genomics



Oncology



Reproductive Health



Genomics in Drug Development

Agrigenomics Research

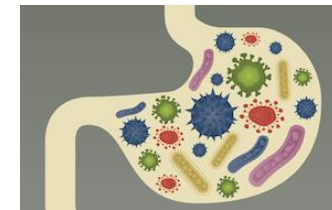


Agrigenomics

Metagenomics Research



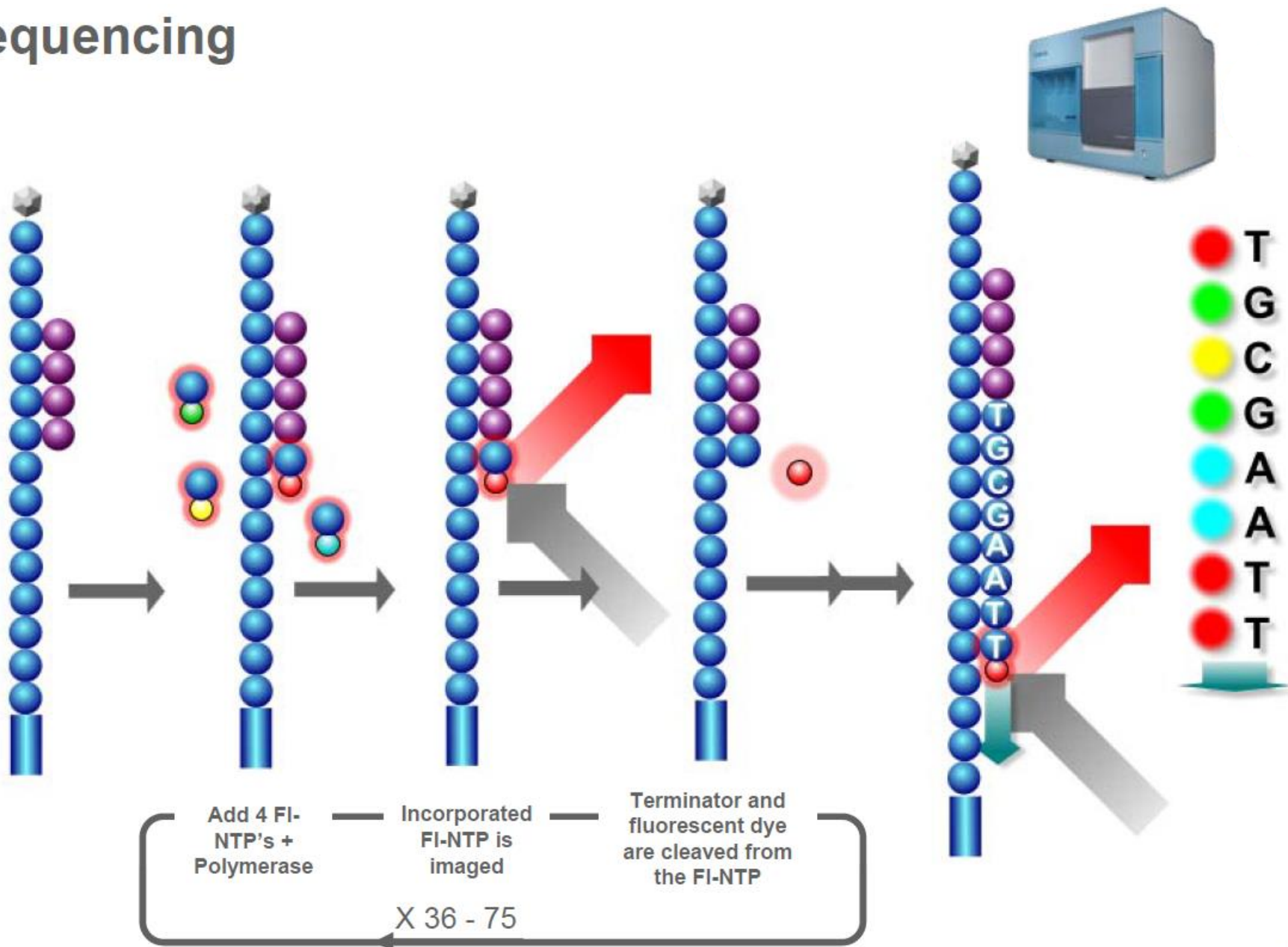
Microbial Genomics



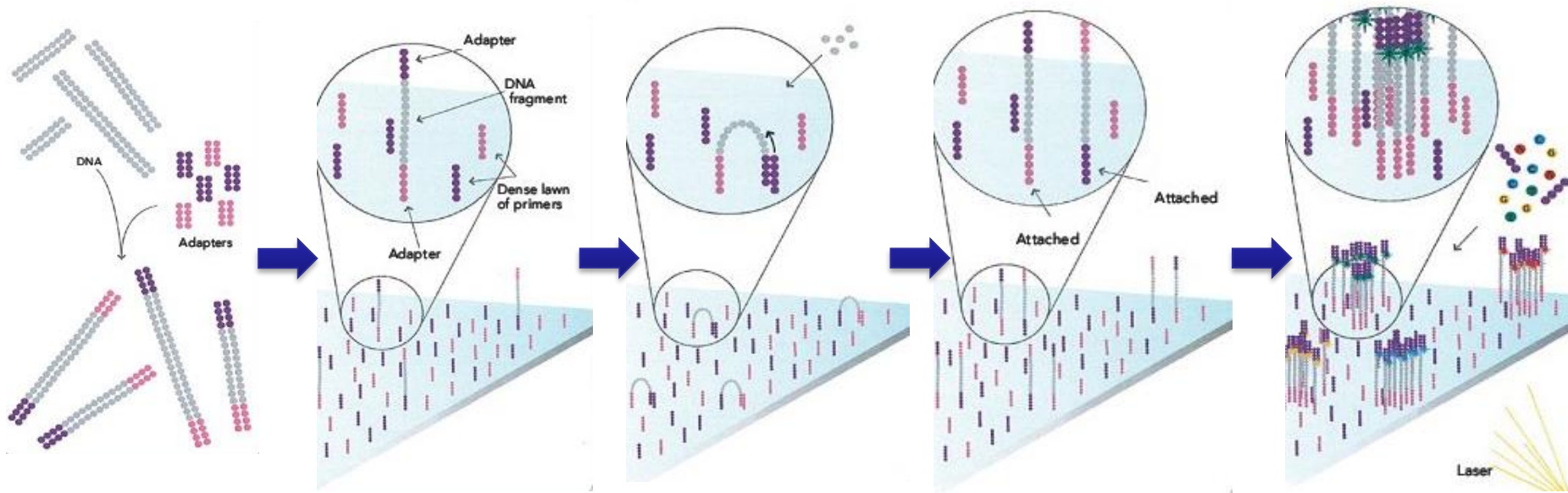
Microbiota



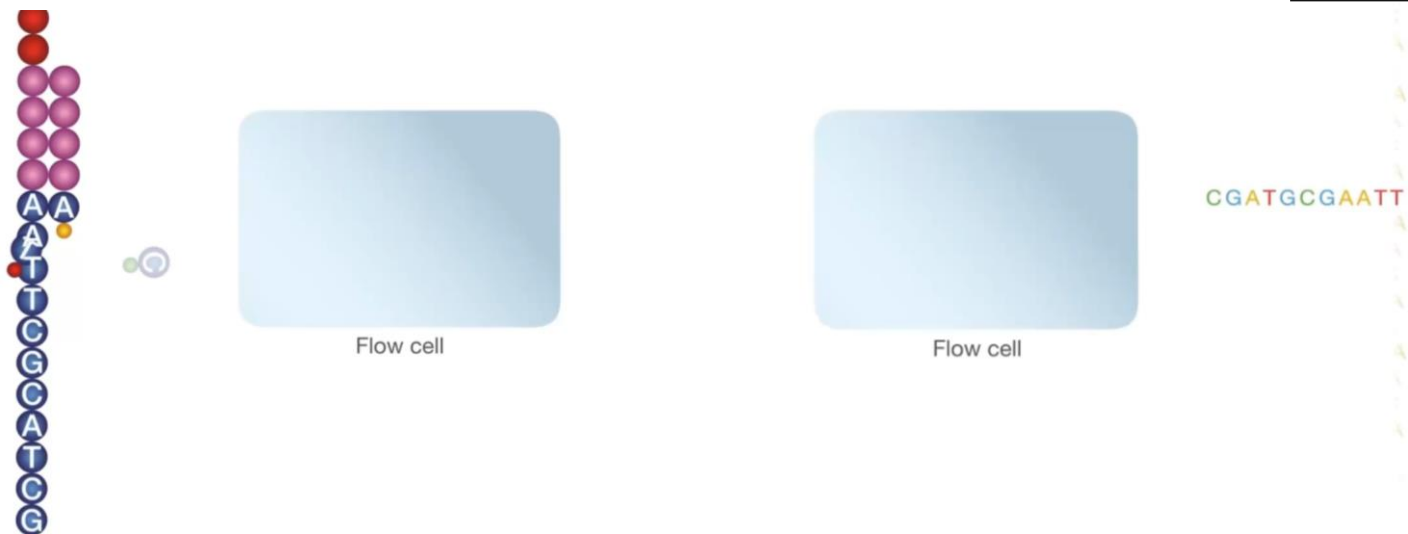
Sequencing



Next-generation sequencing (NGS)

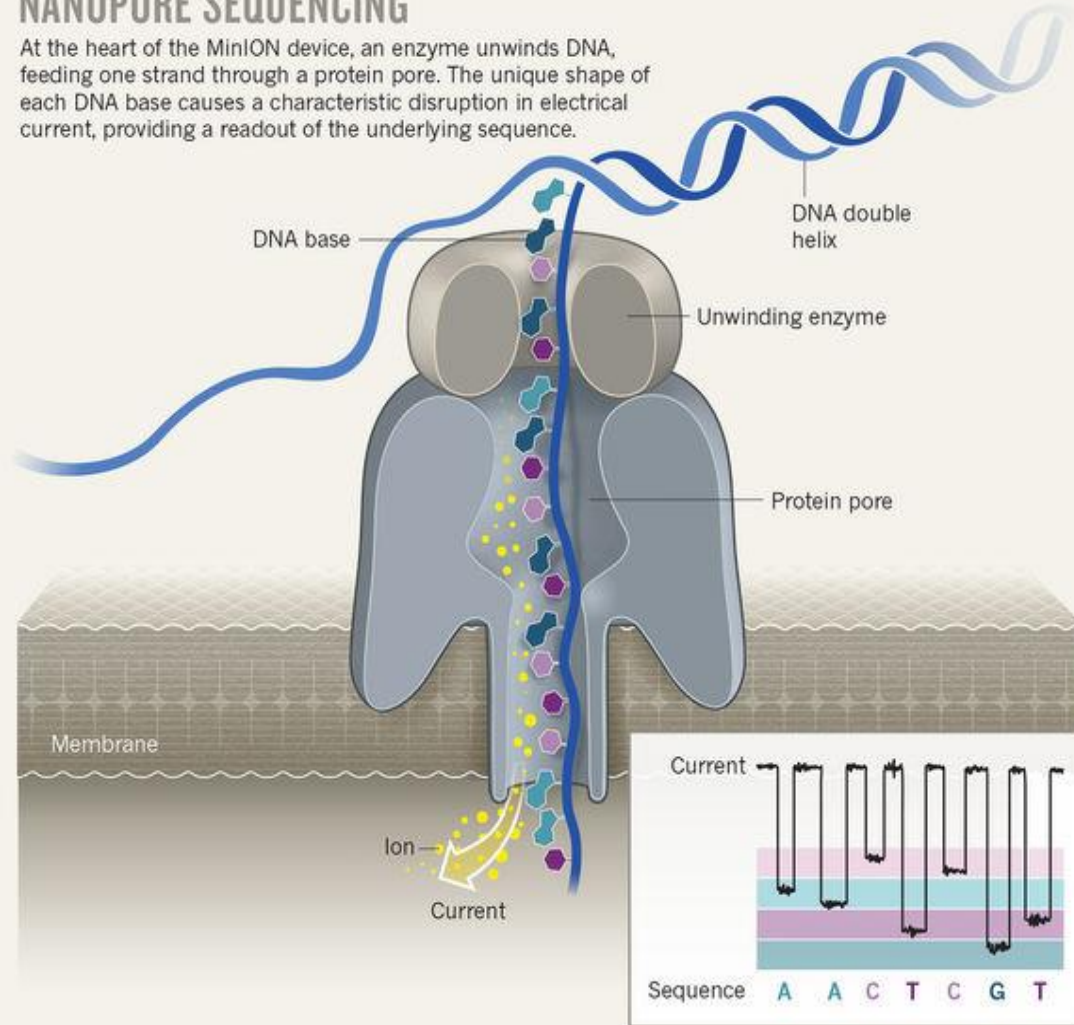


From www.illumina.com



NANOPORE SEQUENCING

At the heart of the MinION device, an enzyme unwinds DNA, feeding one strand through a protein pore. The unique shape of each DNA base causes a characteristic disruption in electrical current, providing a readout of the underlying sequence.



Oxford Nanopore Technology (ONT)



- 超長讀長，取長補短
- 隨測隨停，通量靈活
- 鹼基修飾，實時讀取
- RNA分子，直接測序

Assessing the quality of a Sequencing Read

Phred quality scores

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90 %
20	1 in 100	99 %
30	1 in 1000	99.9 %
40	1 in 10000	99.99 %
50	1 in 100000	99.999 %

$$Q = -10 \log_{10}(P)$$

Q: Phred quality score
P: Base Call error rate

Read of fastq format

```
@SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
GGGTGATGGCCGCTGCCGATGGCGTCAAATCCCACC
+SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII9IG9IC
```



Phred Q score

Sequencing Platform Comparison

Platform Comparison	Illumina	Thermo Fisher	Pacific Biosciences	Oxford Nanopore
Sequencing by Synthesis	Yes	Yes	Yes	No
DNA Size Selection	Yes	Yes	Yes	No
Post-library Amplification	Yes	Yes	No, Single Molecule	No, Single Molecule
Detection	Fluorescent Imaging	Ion Semiconductor	Fluorescent Imaging	Ionic Current Change
Sequencing Rate (s/base)	2 – 20 sec	30 sec	0.25 sec	0.002 sec
Running Time	Fixed	Fixed	Fixed	Run & Stop
DNA Sequencing	Yes	Yes	Yes	Yes
Direct DNA Modification Detection	No	No	Yes	Yes
Direct RNA Sequencing	No	No	No	Yes
Read Length	Short, up to 300 bp X 2	Short, up to 600 bp	Long, Average 6-8 Kb	Long, Average 6-30 Kb
Total Reads (M)	4 – 800 (PE)	2 – 130	0.3 – 0.5	0.3 – 0.5
Total Base (Gb)	1.2 - 120	0.3 – 25 /Chip	5 – 8 /SMRT Cell	2 – 10 /Flow Cell
Instrument Cost (USD)	20K – 275K	200 – 300K	350K	1K – 125K

What Short Read NGS Cannot Do?

LOW DIAGNOSTIC YIELD OF CURRENT STATE-OF-THE-ART NGS BASED TEST

25%

Jamuar and Tan *Human Genomics* (2015) 9:10
DOI 10.1186/s40246-015-0031-5

REVIEW **Open Access**

Clinical application of next-generation sequencing for Mendelian diseases

Saumya Shekhar Jamuar^{1,2} and Ene-Choo Tan^{3,4*}

Abstract

Over the past decade, next-generation sequencing (NGS) has led to an exponential increase in our understanding of the genetic basis of Mendelian diseases. NGS allows for the analysis of multiple regions of the genome in one single reaction and has been shown to be a cost-effective and efficient tool in investigating patients with Mendelian diseases. **More recently, NGS has been successfully deployed in the clinics, with a reported diagnostic yield of ~25%.** However, recommendations on clinical implementation of NGS are still evolving with numerous key challenges that impede the widespread use of genetics in everyday medicine. These challenges include when to order, on whom to order, what type of test to order, and how to interpret and communicate the results, including incidental findings, to the patient and family. In this review, we discuss these challenges and suggest guidelines on implementing NGS in the routine clinical workflow.

Diagnostic yield: The likelihood that a test or procedure will provide the information needed to establish a diagnosis.

Hum Genomics. 2015; 9(1): 10.



complex structural variations

STRUCTURAL VARIATION MORE IMPORTANT THAN SNP

EXPERT OPINION ON DRUG METABOLISM & TOXICOLOGY, 2016
VOL. 12, NO. 2, 135-147
<http://dx.doi.org/10.1517/17425255.2016.1133586>

Taylor & Francis
Taylor & Francis Group

REVIEW

Structural variants can be more informative for disease diagnostics, prognostics and translation than current SNP mapping and exon sequencing

Allen D. Roses^{ab}, P. Anthony Akkari^c, Ornit Chiba-Falek^d, Michael W. Lutz^d, William Kirby Gottschalk^d, Ann Marie Saunders^e, Bob Saul^g, Scott Sundseth^f and Daniel Burns^g

^aDepartment of Neurology and Neurosciences, Duke University, Durham, NC, USA; ^bZinzel Pharmaceuticals, Chapel Hill, NC, USA; ^cShiraz Pharmaceuticals, Inc, Chapel Hill, NC, USA; ^dDepartment of Neurology, Duke University, Durham, NC, USA; ^ePolymorphic DNA, Alameda, CA, USA; ^fCaberner Pharmaceuticals, Inc, Chapel Hill, NC, USA; ^gZinzel Pharmaceuticals, Inc, Raleigh-Durham, NC, USA

ABSTRACT

Introduction: In this article we discuss several human neurological diseases and their relationship to specific highly polymorphic small structural variants (SVs). Unlike genome-wide association analysis (GWAS), this methodology is not a genome screen to define new possibly associated genes, requiring statistical corrections for a million association tests. SVs provide local mapping information at a specific locus. Used with phylogenetic analysis, the specific association of length variants can be mapped and recognized.

Areas covered: This experimental strategy provides identification of DNA variants, particularly variable length Simple Sequence Repeats (SSRs) or STRs or microsatellites) that provide specific local association data at the SV locus. Phylogenetic analysis that includes the specific appearance

ARTICLE HISTORY

Received 31 July 2015
Accepted 14 December 2015
Published online
2 February 2016

KEYWORDS

Alzheimer's disease;
amyotrophic lateral sclerosis;
Lewy Bodies; mitochondrial
metabolism; structural

Expert Opin Drug Metab Toxicol. 2016;12(2):135-47.

Reference	1	2	3			
Insertion	1	2	5	3		
Deletion	1	3				
Inversion	1	3	2			
Copy Number Variation	1	1	1	1	2	3
Tandem Duplication	1	1	2	3		
Dispersed Duplication	1	2	1	3		
Mobile Element Insertion	1	2	Mobile Element	3		
Translocation	1					
	10	11	12	2	3	

https://en.wikipedia.org/wiki/Human_Genome_Structural_Variation

Structural Variants are More Impactful than SNP

EXPERT OPINION ON DRUG METABOLISM & TOXICOLOGY, 2016
 VOL. 12, NO. 2, 135-147
<http://dx.doi.org/10.1517/17425255.2016.1133586>



REVIEW

Structural variants can be more informative for disease diagnostics, prognostics and translation than current SNP mapping and exon sequencing

Allen D. Roses^{ab}, P. Anthony Akkari^c, Ornit C. Ann Marie Saunders^d, Bob Saul^e, Scott Sund

^aDepartment of Neurology and Neurosciences, Duke University Medical Center, Durham, NC, USA; ^bDepartment of Neurology, Duke University Medical Center, Durham, NC, USA; ^cCaberner Pharmaceuticals, Inc, Chapel Hill, NC, USA; ^dDepartment of Neurology, Duke University Medical Center, Durham, NC, USA; ^eDepartment of Neurology, Duke University Medical Center, Durham, NC, USA

ABSTRACT

Introduction: In this article we discuss several highly polymorphic small structural variants (SVs) that are highly polymorphic in GWAS, this methodology is not a germline variant, requiring statistical corrections for a million variants can be mapped and recognized.

Areas covered: This experimental strategy provides variable length Simple Sequence Repeats (SSRs) and local association data at the SV locus. *Pharmacogenomics*

Expert Opin Drug Metab Toxicol. 2016;12(2):135-47.

ORIGINS OF REDUCED ACCURACY IN CLINICAL GENOMICS FROM SHORT SEQUENCING READS



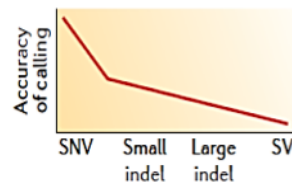
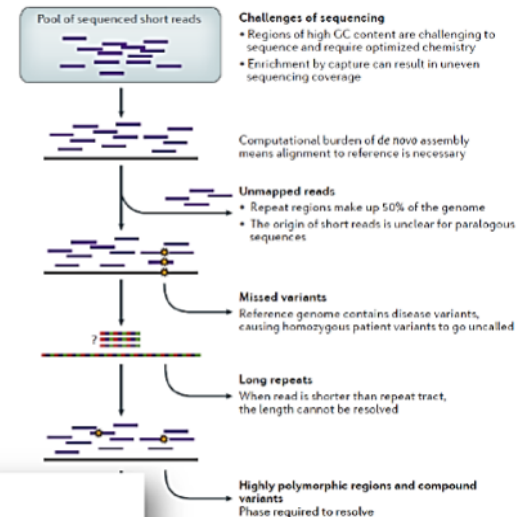
REVIEWS

APPLICATIONS OF NEXT-GENERATION SEQUENCING

Towards precision medicine

Euan A. Ashley

Abstract | There is great potential for genome sequencing to enhance patient care through improved diagnostic sensitivity and more precise therapeutic targeting. To maximize this potential, genomics strategies that have been developed for genetic discovery—including DNA sequencing technologies and analysis algorithms—need to be adapted to fit clinical needs. This will require the optimization of alignment algorithms, attention to quality-coverage metrics, tailored solutions for paralogous or low-complexity areas of the genome, and the adoption of consensus standards for variant calling and interpretation. Global sharing of this more accurate genotypic and phenotypic data will accelerate the determination of causality for novel genes or variants. Thus, a deeper understanding of disease will be realized that will allow its targeting with much greater therapeutic precision.



Accuracy of variant calling falls with increasing disruption of the open reading frame

Nat Rev Genet. 2016 Aug 16;17(9):507-22.

We Need a Deeper Genome

LOW *DIAGNOSTIC YIELD OF CURRENT STATE-OF-THE-ART NGS BASED TEST

Jamuar and Tan *Human Genomics* (2015) 9:10
DOI 10.1186/s40246-015-0031-5



Human Genomics

REVIEW

Open Access



Clinical application of next-generation sequencing for Mendelian diseases

Saumya Shekhar Jamuar^{1,2} and Ene-Choo Tan^{2,3*}

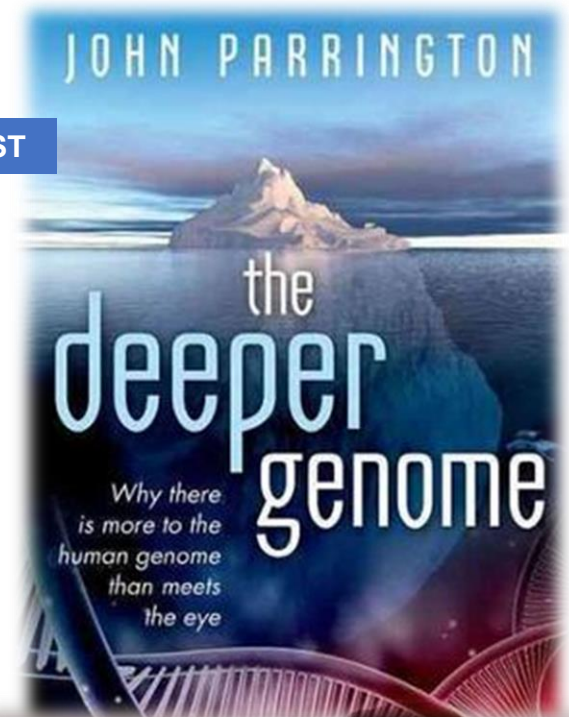
Abstract

Over the past decade, next-generation sequencing (NGS) has led to an exponential increase in our understanding of the genetic basis of Mendelian diseases. NGS allows for the analysis of multiple regions of the genome in one single reaction and has been shown to be a cost-effective and efficient tool in investigating patients with Mendelian diseases. More recently, NGS has been successfully deployed in the clinics, with a reported diagnostic yield of ~25%. However, recommendations on clinical implementation of NGS are still evolving with numerous key challenges that impede the widespread use of genetics in everyday medicine. These challenges include when to order, on whom to order, what type of test to order, and how to interpret and communicate the results, including incidental findings, to the patient and family. In this review, we discuss these challenges and suggest guidelines on implementing NGS in the routine clinical workflow.

Keywords: Next-generation sequencing, Whole exome sequencing, Clinical applications, Mendelian diseases

***Diagnostic Yield:** the likelihood that a test or procedure will provide the information needed to establish a diagnosis

Hum Genomics. 2015; 9(1): 10.



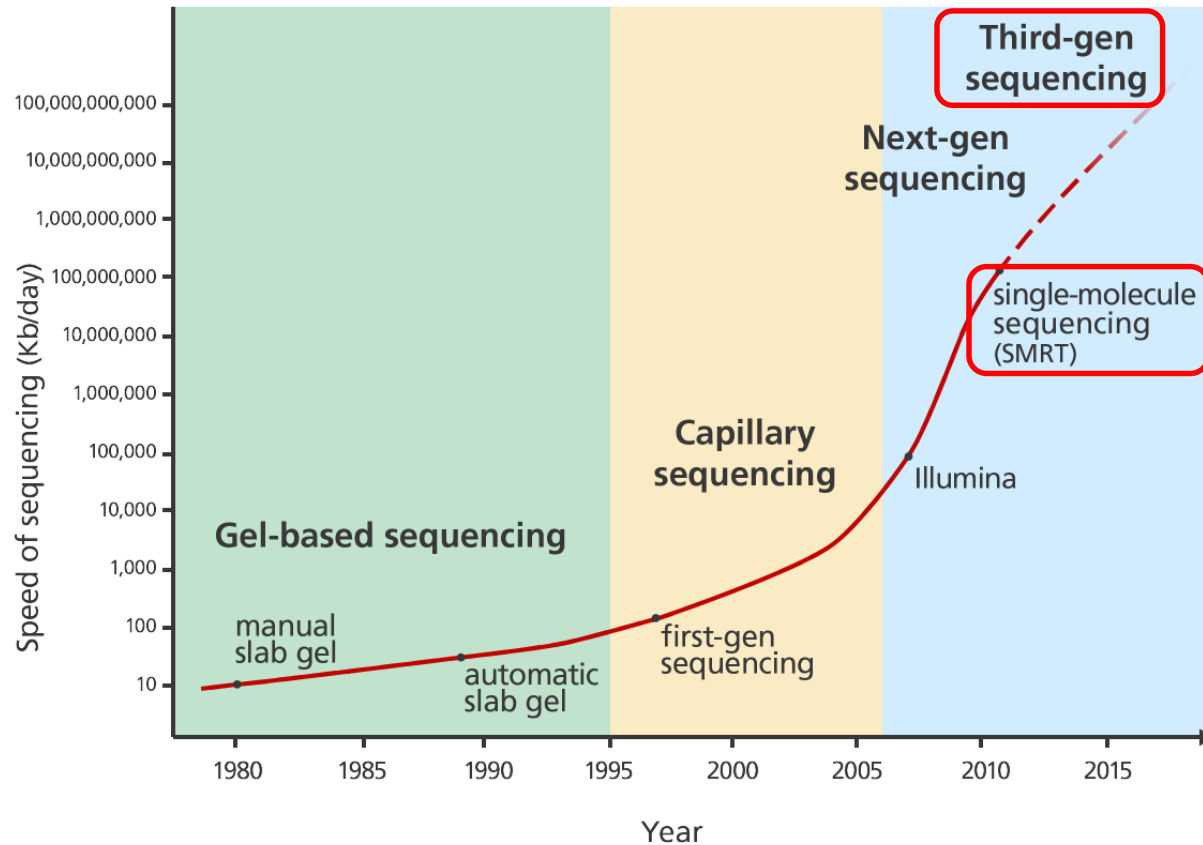
John Parrington

Author of *The Deeper Genome*

https://www.youtube.com/watch?v=WYvuAtQo-_g

What's Beyond Short Read NGS?

HUMAN GENOMIC REGIONS UNRESOLVED BY SHORT SEQUENCING READS



<https://www.yourgenome.org/stories/third-generation-sequencing>

- Repeated regions (simple repeats, tandem repeats, transposon-related repeats)
- Highly polymorphic regions (HLA) – haplotype phasing
- Structural variants (relocation, inversion, duplications)
- Un-sequencable regions by NGS (AT or GC rich regions)

How “Long Reads” can Help?

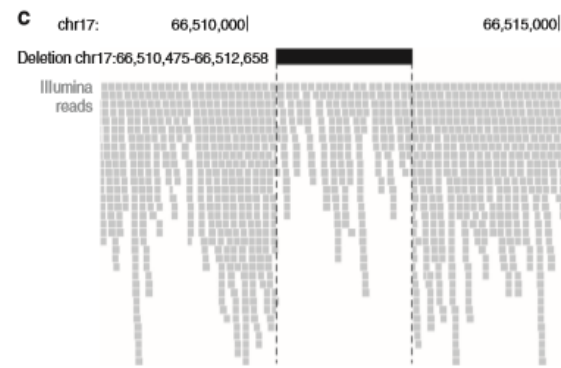
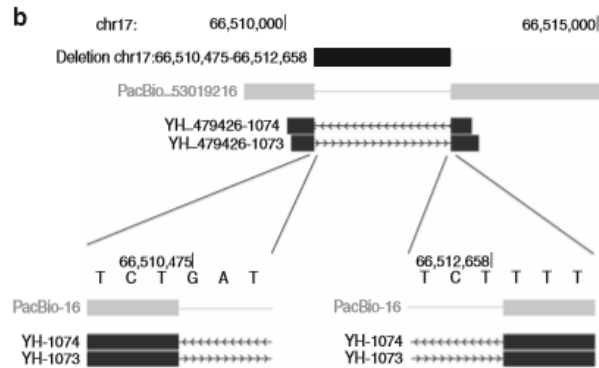
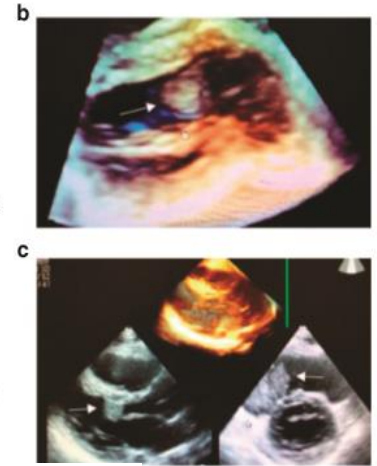
© American College of Medical Genetics and Genomics

BRIEF REPORT | **Genetics in Medicine**

Long-read genome sequencing identifies causal structural variation in a Mendelian disease

Jason D. Merker, MD, PhD^{1,2}, Aaron M. Wenger, PhD³, Tam Sneddon, DPhil², Megan Grove, MS, LCGC², Zachary Zappala, PhD^{1,4}, Laure Fresard, PhD¹, Daryl Waggott, MSc^{5,6}, Sowmi Utiramerur, MS², Yanli Hou, PhD¹, Kevin S. Smith, PhD¹, Stephen B. Montgomery, PhD^{1,4}, Matthew Wheeler, MD, PhD^{5,6}, Jillian G. Buchan, PhD^{1,2}, Christine C. Lambert, BA³, Kevin S. Eng, MS³, Luke Hickey, BS³, Jonas Korlach, PhD³, James Ford, MD^{4,5,7} and Euan A. Ashley, MRCP, DPhil^{2,4,5,6}

- a**
- 7 yrs ● Left atrial myxoma resection, atrial repair
 - 10 yrs ● Testicular mass, right orchiectomy
 - 13 yrs ● Pituitary tumour
 - 16 yrs ● Recurrence of myxomata, resection, adrenal microadenoma
 - 18 yrs ● Recurrence of ventricular myxomata, resection, VT
 - 19 yrs ● ACTH-independent Cushing's disease, thyroid nodules
 - 21 yrs ● Transsphenoidal resection of pituitary
 - Present ● Recurrence of myxomata, under consideration for heart transplant



Features and Limitations of Nanopore Sequencing

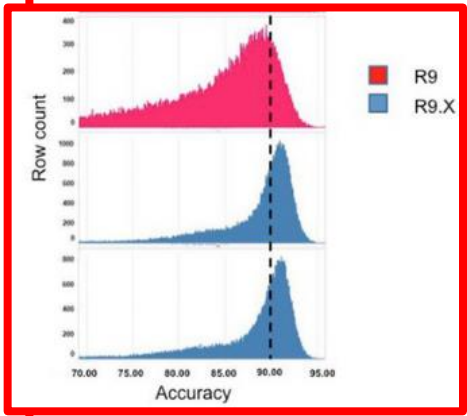
- Very long read-length
- Low capital investment
- Single molecule, PCR-free
- Real-time & urgent applications

E. coli assembly in 8 reads



Joshua Quick,
University of Birmingham,
Feb 15, 2018

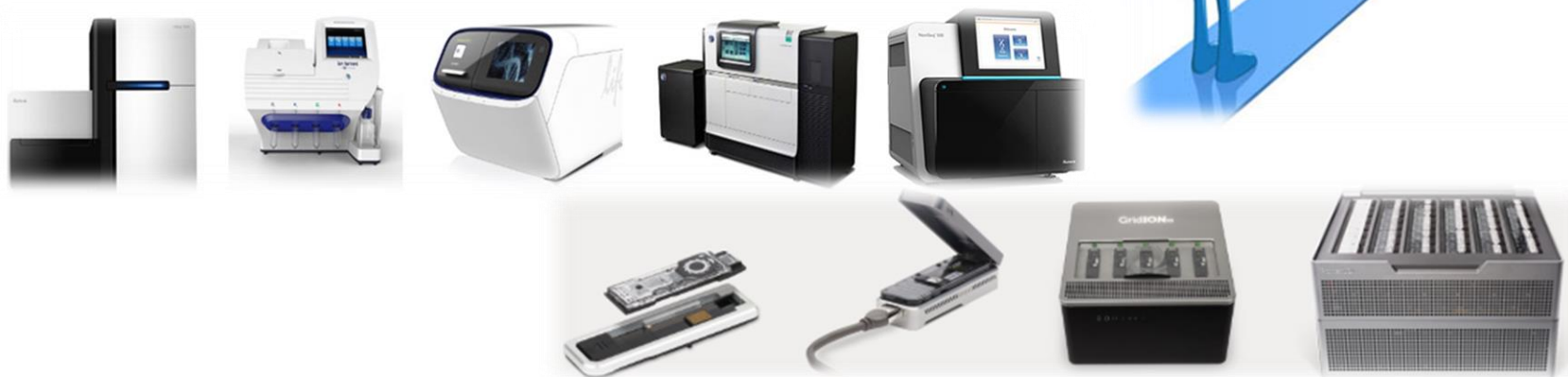
- Lower throughput
- Higher cost per GB than NGS
- Lower sequencing quality than NGS
- Hard calling genetic variants and mutations



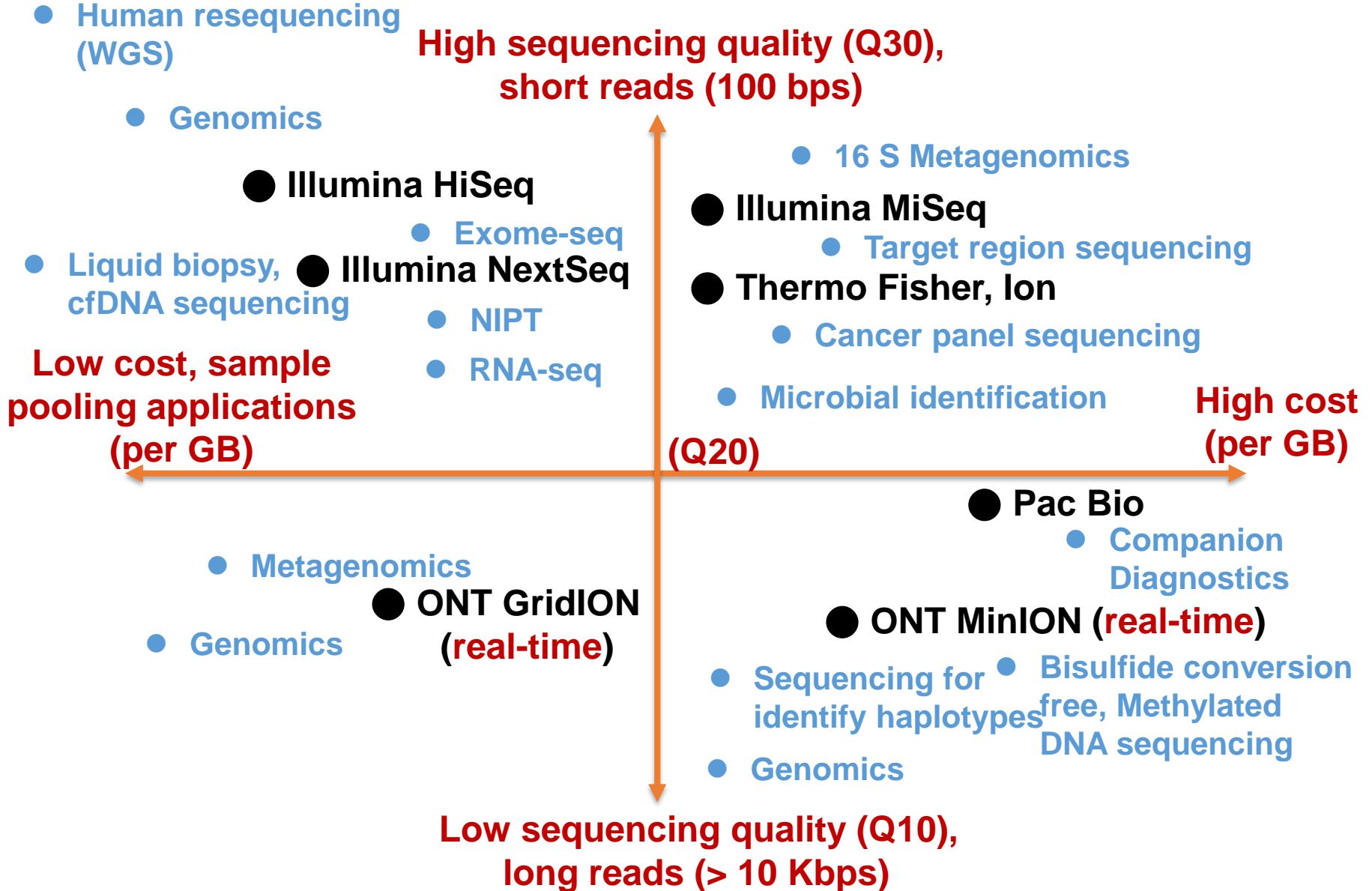
Clive G. Brown, 2016

QUESTION!

- For medical applications
- For academic research applications
- Key questions
 - How to choose platforms?
 - Will it be appropriate to integrated two platforms for an application?



Strategy for Selecting Seq. Platforms



Application-1

Exploring Human Genome Structural Variants (SV)

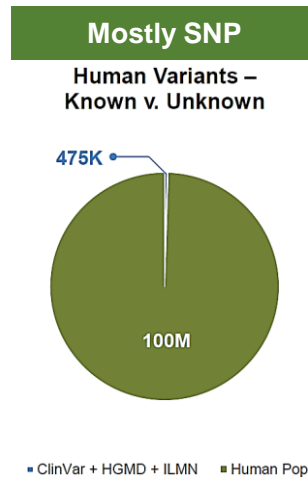
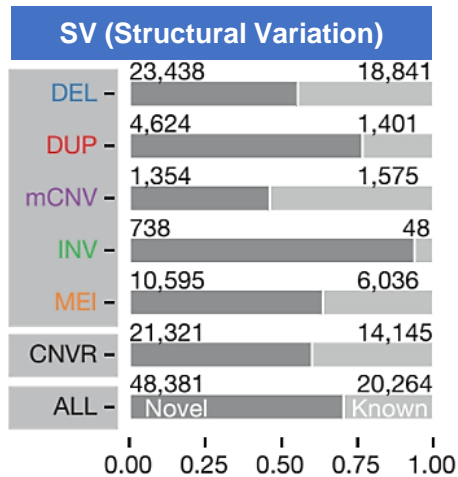
Human SV – Current Knowledge

ARTICLE

OPEN
doi:10.1038/nature15394

An integrated map of structural variation in 2,504 human genomes

A list of authors and their affiliations appears at the end of the paper.



~ 4 K
~ 18.4 Mbp

Individual Level: ~ 4 M
~ 4 Mbp

Nature. 2015 Oct 1;526(7571):75-81.

Catalog of Human SV:

- Illumina WGS (~100 bp reads, mean 7.4X coverage)
- 26 human populations

“...although SNPs contribute more eQTLs overall, our results suggest that SVs have a disproportionate impact on gene expression relative to their number.”

“We further expand the number of candidate SVs in strong LD with GWAS hits by ~30% (39/136 novel associations implicating SVs as candidates) and find that GWAS haplotypes are enriched up to threefold for common SVs, which emphasizes the relevance of ascertaining SVs in disease studies.”

“...while many SVs in our callset are statistically phased, the diploid nature of the genome is non-optimally captured by current analysis approaches, which mostly rely on mapping to a haploid reference. We envision that in the future, the use of technology allowing substantial increases in read lengths over the current state-of-the-art will enable genomic analyses of truly diploid sequences to facilitate targeting these additional layers of genomic complexity.”

Human Genome Assembly by ONT Reads

ARTICLES

nature
biotechnology

OPEN

Nanopore sequencing and assembly of a human genome with ultra-long reads

Miten Jain^{1,13}, Sergey Koren^{2,13}, Karen H Miga^{1,13}, Josh Quick^{3,13}, Arthur C Rand^{1,13}, Thomas A Sasaki^{4,5,13}, John R Tyson^{6,13}, Andrew D Beggs⁷, Alexander T Dilthey², Ian T Fiddes¹, Sunir Malla⁸, Hannah Marriott⁸, Tom Nieto⁷, Justin O'Grady⁹, Hugh E Olsen¹, Brent S Pedersen^{4,5}, Arang Rhie², Hollian Richardson⁹, Aaron R Quinlan^{4,5,10}, Terrance P Snutch⁶, Louise Tee⁷, Benedict Paten¹, Adam M Phillippy², Jared T Simpson^{11,12}, Nicholas J Loman³ & Matthew Loose⁸

We report the sequencing and assembly of a reference genome for the human GM12878 Utah/Ceph cell line using the MinION (Oxford Nanopore Technologies) nanopore sequencer. 91.2 Gb of sequence data, representing ~30x theoretical coverage, were produced. Reference-based alignment enabled detection of large structural variants and epigenetic modifications. *De novo* assembly of nanopore reads alone yielded a contiguous assembly (NG50 ~3 Mb). We developed a protocol to generate ultra-long reads (N50 > 100 kb, read lengths up to 882 kb). Incorporating an additional 5x coverage of these ultra-long reads more than doubled the assembly contiguity (NG50 ~6.4 Mb). The final assembled genome was 2,867 million bases in size, covering 85.8% of the reference. Assembly accuracy, after incorporating complementary short-read sequencing data, exceeded 99.8%. Ultra-long reads enabled assembly and phasing of the 4-Mb major histocompatibility complex (MHC) locus in its entirety, measurement of telomere repeat length, and closure of gaps in the reference human genome assembly GRCh38.

Nat Biotechnol. 2018 Apr;36(4):338-345.

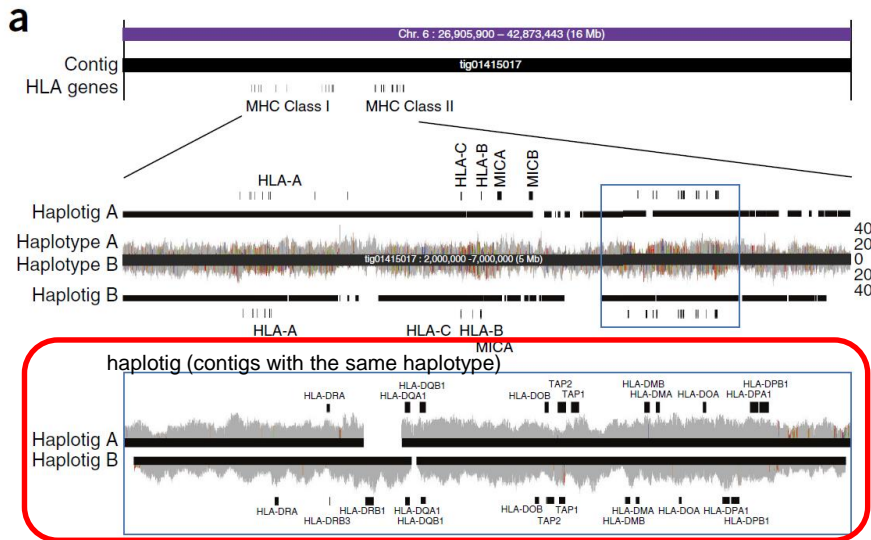
The Final Assembled Genome:

- 41 FC (Regular), 5 FC (Ultra-long)
- 91.2 Gb, ~30X
- Ultra-long Reads, N50 > 100 Kb, up to 882 Kb
- 2,867 million bases
- Covering 85.8% of reference genome
- Assembly accuracy: 99.8% (incorporating complementary short-read data)
- NG50 ~3 Mb
- With Ultra-long reads, NG50 ~6.4 Mb
- With Ultra-long reads, achieved phasing of the entire 4 Mb MHC locus

Ultra-long Reads, Assembly and Telomeres

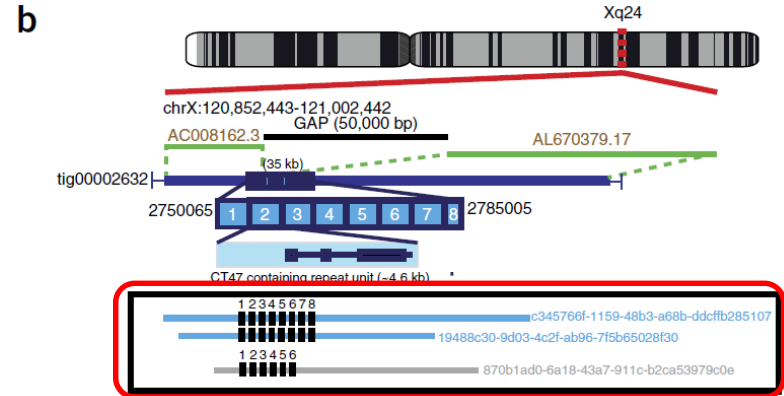
nature
biotechnology

- OPEN • Heterozygous SNPs were called using Illumina data
 • Phased using the ultra-long nanopore reads
 • Generate two **pseudo-haplotypes**

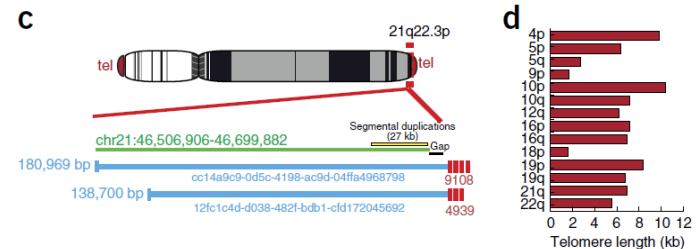


- A 16-Mbp ultra-long read contig and associated haplotigs are shown spanning the full MHC region.
- The first time the MHC has been assembled and phased over its full length in a diploid human genome.

Nat Biotechnol. 2018 Apr;36(4):338-345.



Two reads provide evidence for an array of eight repeat copies and one read supports six copies, suggesting heterozygosity.



- FISH (fluorescent in situ hybridization) estimates and direct cloning of telomeric DNA suggests that telomere repeats (**TTAGGG**) extend for multiple kbs at the ends of each chromosome.
- Evidence for telomeric arrays that span 2–11 kb within 14 subtelomeric regions for GM12878.

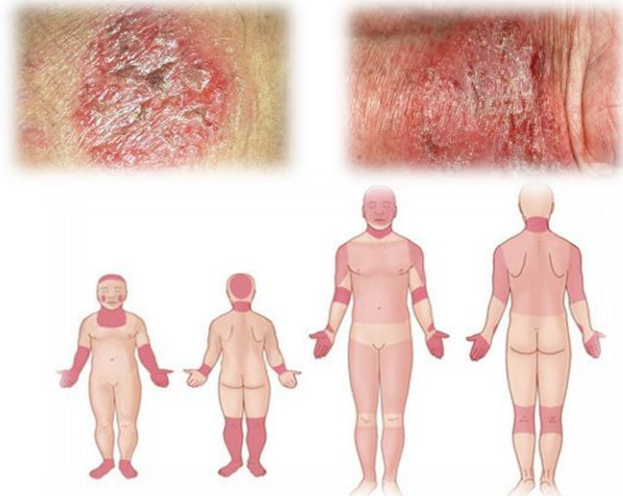
One remarkable molecule: Filaggrin

- Ichthyosis vulgaris (尋常魚鱗病); Eczema (濕疹)
- The CNV allele frequencies in the **Irish** population were found to be 33.9% **10** repeats, 51.5% **11** repeats and 14.6% **12** repeats. Shortest genotype (10,10): Eczema risk: **1.67**
- When null mutations are excluded, each additional filaggrin repeat gained decreases the odds ratio for atopic eczema by **0.88**.
- Filaggrin CNV makes a significant, **dose-dependent** contribution to eczema risk

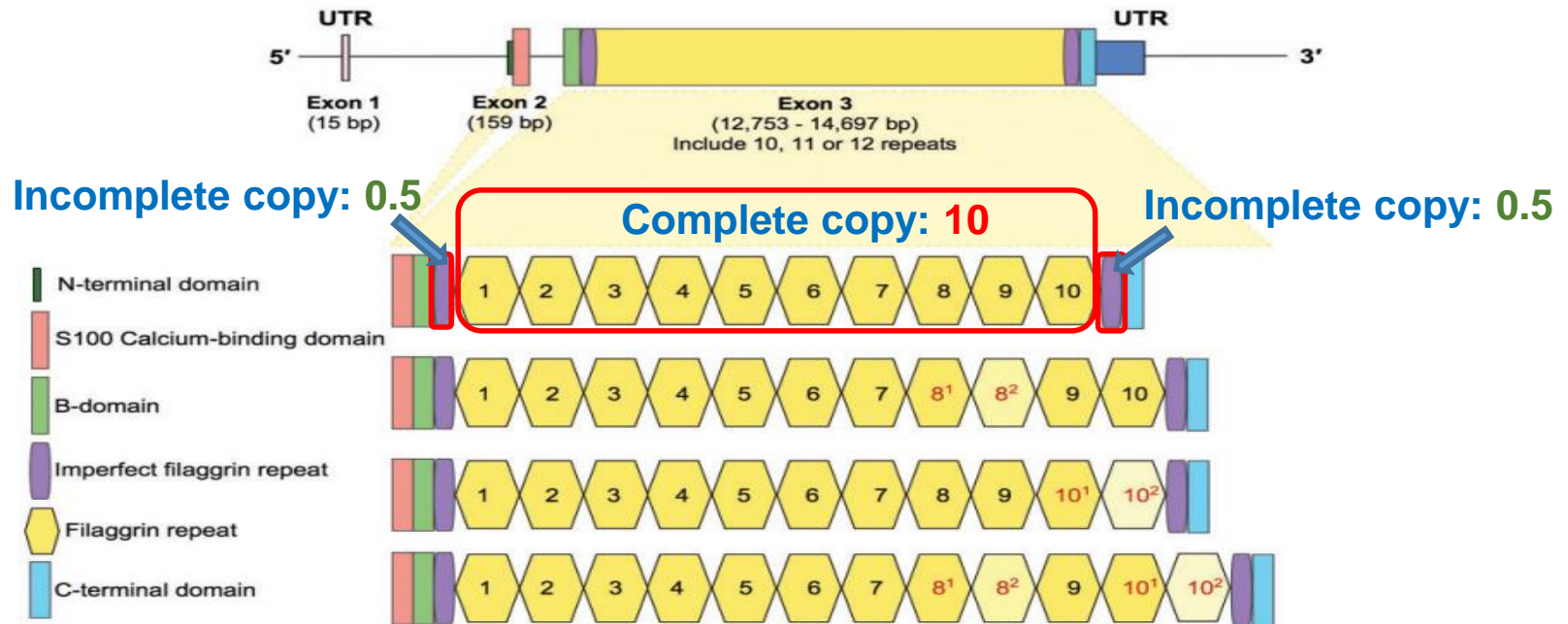
Ichthyosis vulgaris (魚鱗病)



Eczema (濕疹)



Structure of FLG and Copy number



- Filaggrin repeats were shown to be copy number variable ranging from **10 to 12 copies** among human populations.
- The **copy number** of these repeats was **negatively associated** with atopic dermatitis susceptibility.
- However, a comprehensive documentation of the global distribution of FLG genetic variation free of ascertainment bias has yet to be compiled.

Tandem Repeats Variations

COMMENTARY

See related article on pg 98

Profilaggrin, Dry Skin, and Atopic Dermatitis Risk: Size Matters

John A. |

Mutation risk factors. New findings in *FLG* and utility in

Journal of In

Pediatric Dermatology Vol. 34 No. 3 e140–e141, 2017

Intragenic Copy Number Variation in the Filaggrin Gene in Ethiopian Patients with Atopic Dermatitis

Abstract: Genes involving truncation and copy number variation are associated with the development of atopic dermatitis in Asian populations. Identification of a relationship between *FLG* copy number and atopic dermatitis severity in a small population is proposed. We studied the association between *FLG* copy number and atopic dermatitis severity in Ethiopian patients, suggesting that factors are of importance to AD.

Research letter

Copy-number variation of the filaggrin gene in Korean patients with atopic dermatitis: what really matters, ‘number’ or ‘variation’?

DOI: 10.1111/bjd.14287

DEAR EDITOR, Since the articles reporting a methodological breakthrough on the full sequencing of the gene encoding profilaggrin (*FLG*),^{1,2} associations between loss-of-function muta-

***Br J Dermatol.* 2016 May;174(5):1098-100.**

been suggested that the number of *FLG* repeats can relate to a dry skin phenotype; i.e., fewer repeats may lead to less *FLG* protein expression and drier skin (Ginger *et al.*, 2005). What is currently not known, however, is whether there is a relationship between the intragenic copy-

***J Invest Dermatol.* 2012 Jan;132(1):10-1.**

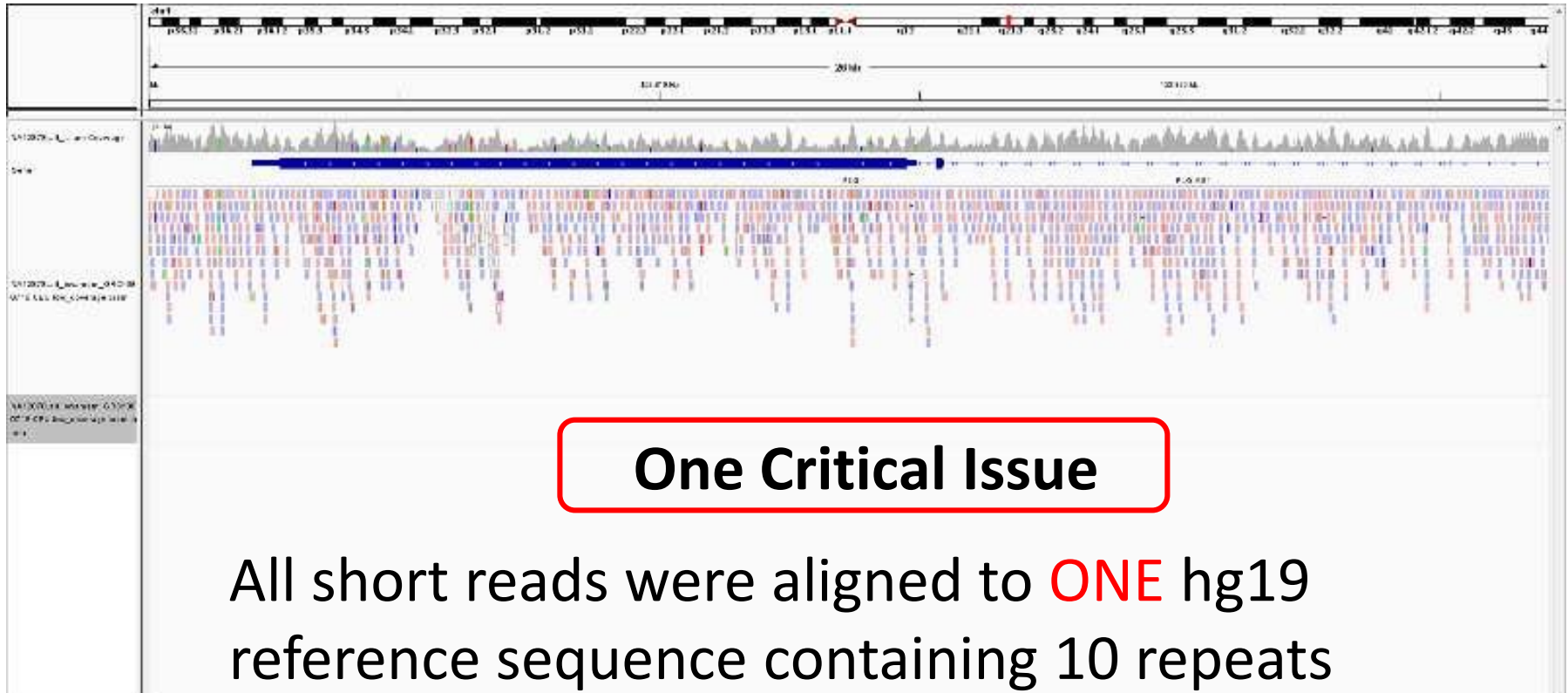
long-range polymerase chain reaction (PCR) amplification of the *FLG* repeats, which has previously been described for various populations, including Africans (4). The total *FLG* CN was classified as low (20–21) or high (22–24).

A total of 105 cases were successfully genotyped and phenotyped and included in the study. The long-range PCR yielded products of sizes that confirmed a different

***Pediatr Dermatol.* 2017 May;34(3):e140-e141.**

tions of *FLG* and atopic dermatitis (AD) have been reported across ethnicities.^{3,4} However, both the low prevalence of *FLG* mutations in patients with AD in some nations (< 4% in Italy) and the high prevalence of *FLG* mutations in healthy control in other nations (~ 10% in Ireland) suggest that factors other than *FLG* mutation may be at work.^{5,6} Brown *et al.* introduced an interesting new factor contributing to the risk of AD: copy-number variation (CNV).⁶ *FLG* is polymorphic, with allelic variants of 10–12 nearly identical repeats in exon 3.⁶ They

Conventional Short Reads Alignment



One Critical Issue

All short reads were aligned to **ONE** hg19 reference sequence containing 10 repeats

Cannot make high confidence variant calls in the repeat region using short reads

Long Range PCR for full length of Filaggrin

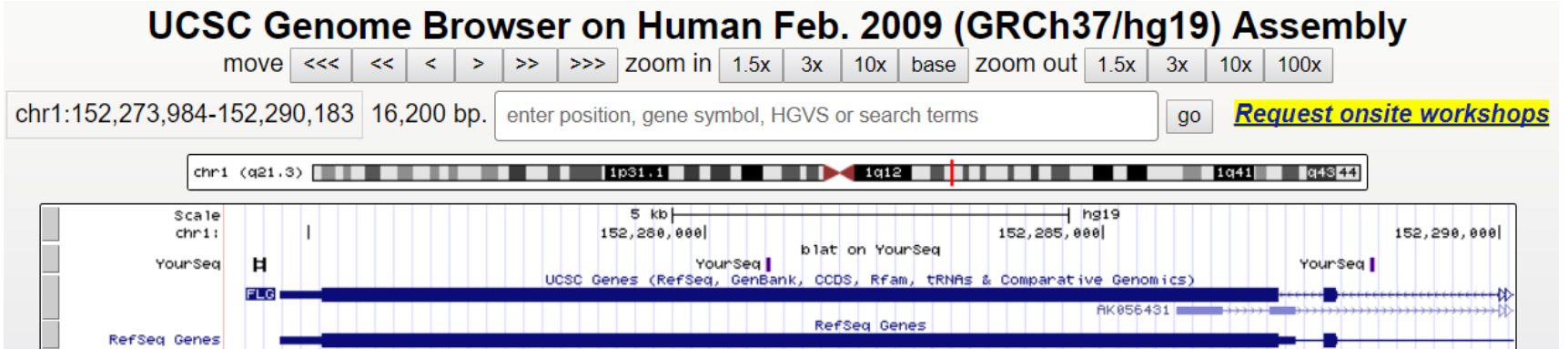
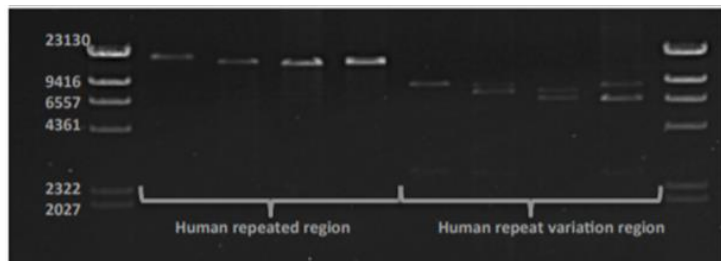
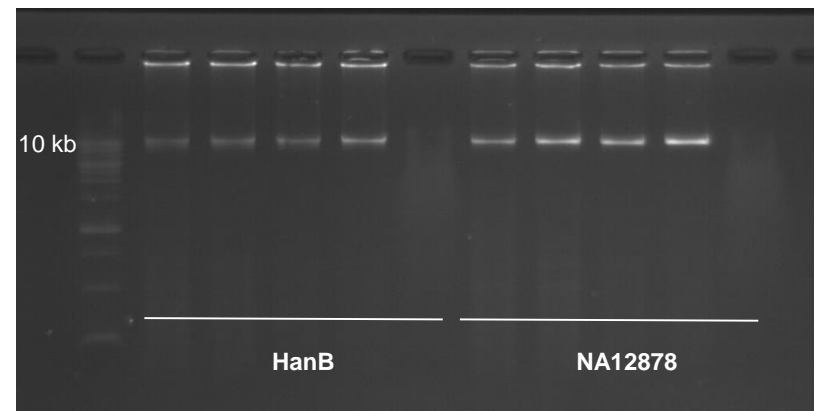


Table S1. Primer sets for the filaggrin gene repeat region in human, chimpanzee, gorilla, orangutan, crab-eating macaque, and human variation in repeated region.

REPEATED REGION	
Human	hFLG (13969 bp in total)
Forward	5' -CTGTGCATATGGCTAACTGGCTTTCAGAGA-3'
Reverse	5' -ATTGTGGGACAGTGATTATGTTGGAGAAAA-3'
Human variation in repeated region	hFLGv (6480 bp in total)
Forward	5' -GTGCAAGCAGAAAAACATATGACA -3'
Reverse	5' -CCTGTTTCGTGATCTGCCTTGACATGG -3'



PCR Amplification

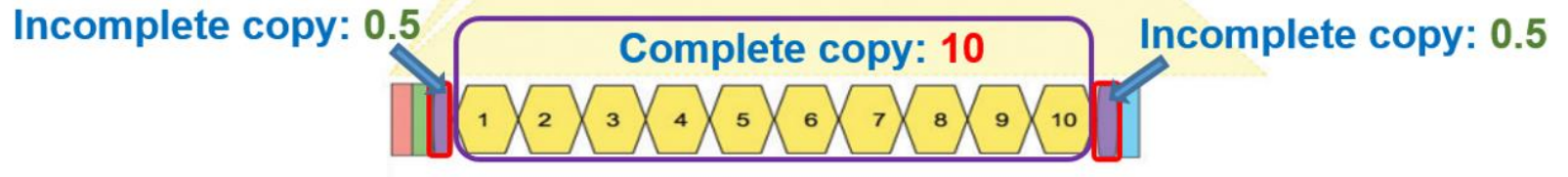


image_IM009993

Preliminary Results

Sample	HanB1	HanB2	HanB3	NA12878
Reads (Q>=7)	3,182	6,246	49,346	277850
Align to FLG ref	4262	5068	34558	9169
>10,000 bp	686	724	3762	2443
Both F and R found	323	445	1780	1533
TRN Determined	(11+1/12+1)	(11+1/12+1) (with noise)	(11+1/12+1)	(10+1/12+1)

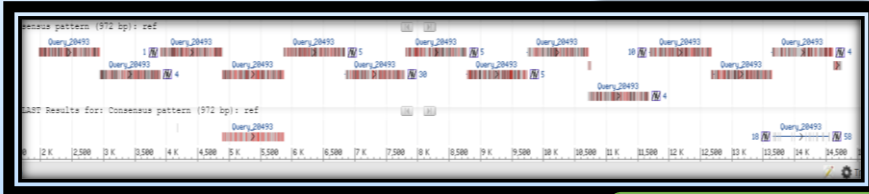
$$10 + 0.5 + 0.5 = 10 + 1$$



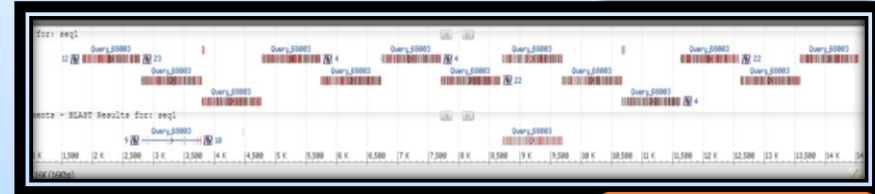
Identifying the copy number of FLG gene

Blast result

12+1 copy



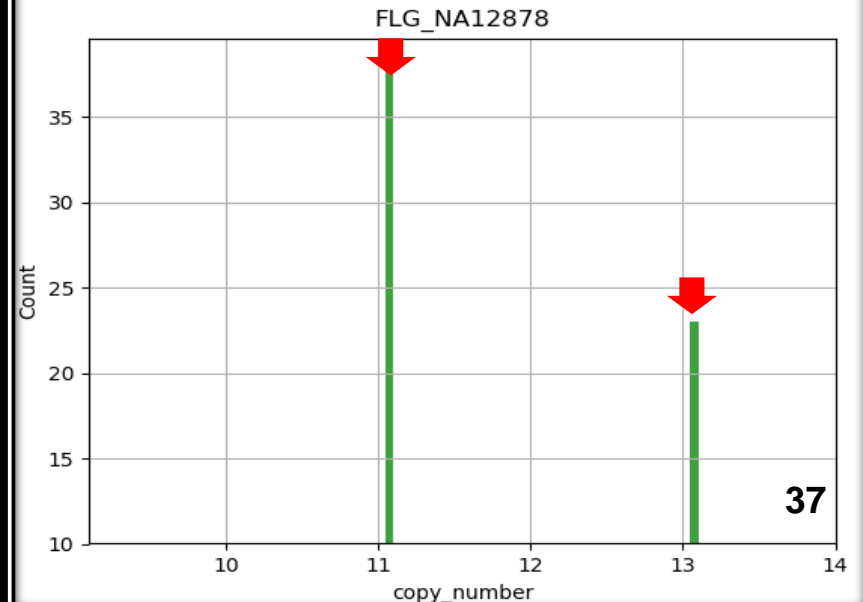
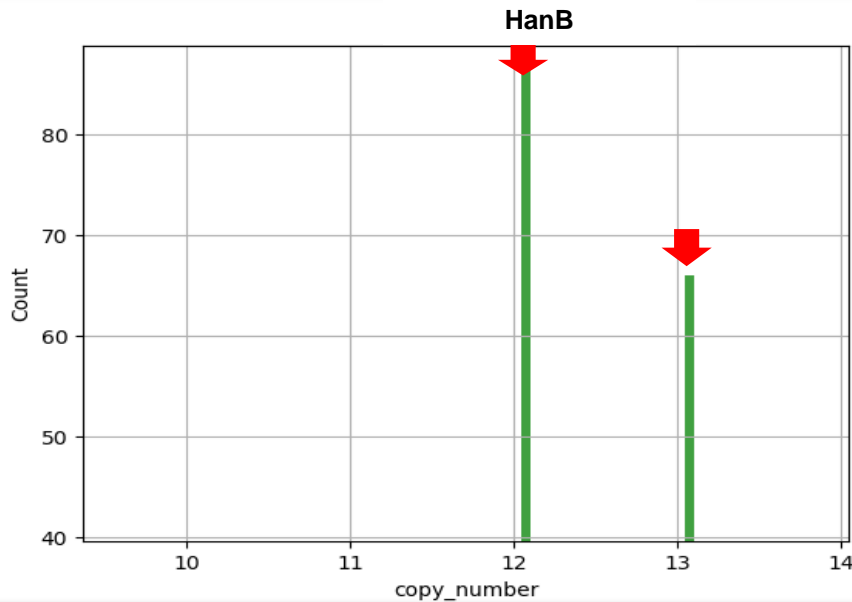
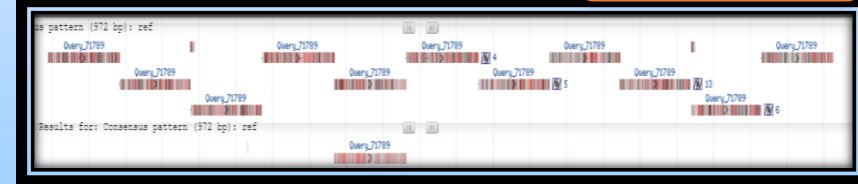
12+1 copy



11+1 copy



10+1 copy

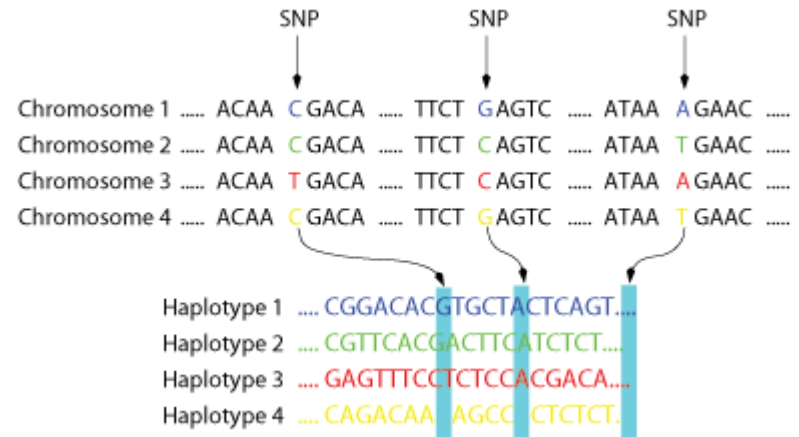
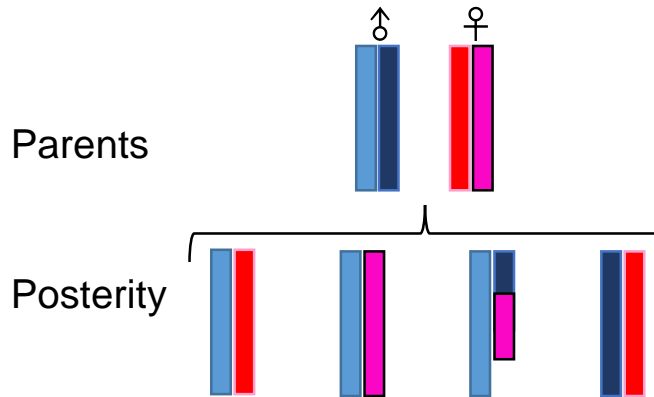


Application-2

Full genotyping for high polymorphic loci in Pharmacogenomics

Haplotypes (Full genotyping of a high polymorphic locus)

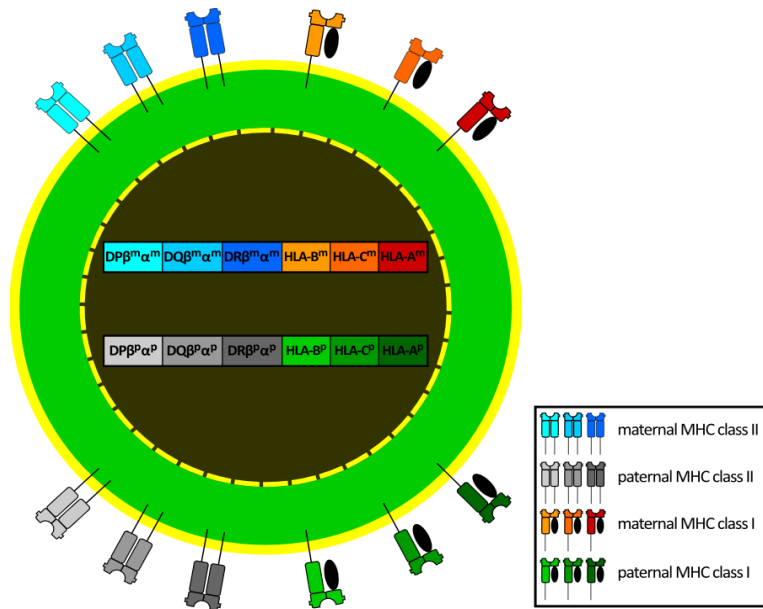
- A haplotype is a group of genes within an organism that was inherited together from a single parent
- The term "**haplotype**" can also refer to the inheritance of a cluster of single nucleotide polymorphisms (SNPs)
- There are highly polymorphic (more haplotypes) on genes of CYP or HLA



Filippo Geraci and Marco Pellegrini

Human Leukocyte Antigen (HLA)

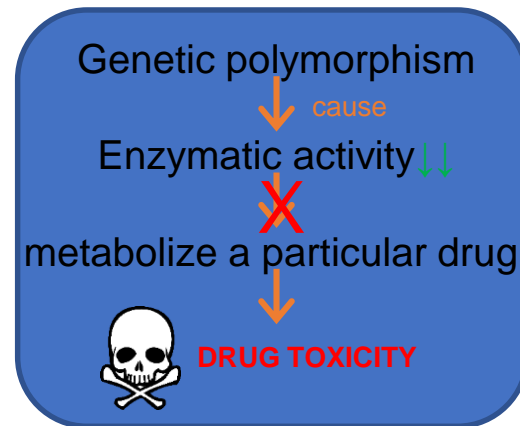
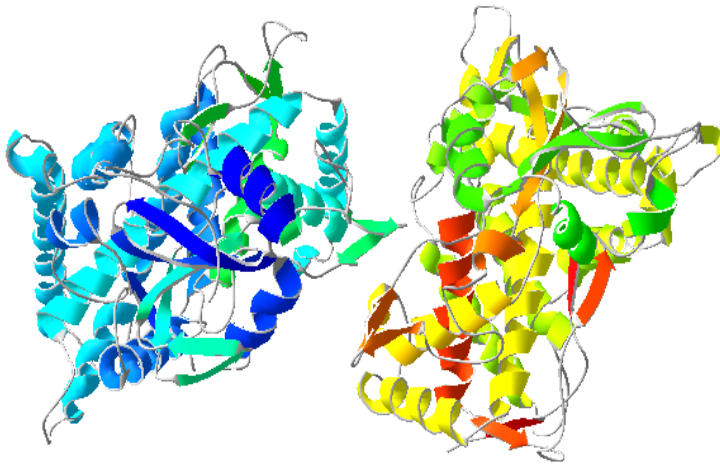
- The **h**uman **l**eukocyte **a**ntigen (HLA) genes are the human versions of the **m**ajor **h**istocompatibility **c**omplex (MHC)
- HLA system is the locus of genes that encode for proteins on the surface of cells that are responsible for regulation of the immune system in humans



Gene	Number of Alleles (Haplotypes)
HLA-A	4,081
HLA-B	4,950
HLA-C	3,685
HLA-DQA1	94
HLA-DQB1	1,178
HLA-DPA1	64
HLA-DPB1	963
HLA-DRB1	2,146

Cytochrome P450 Enzymes

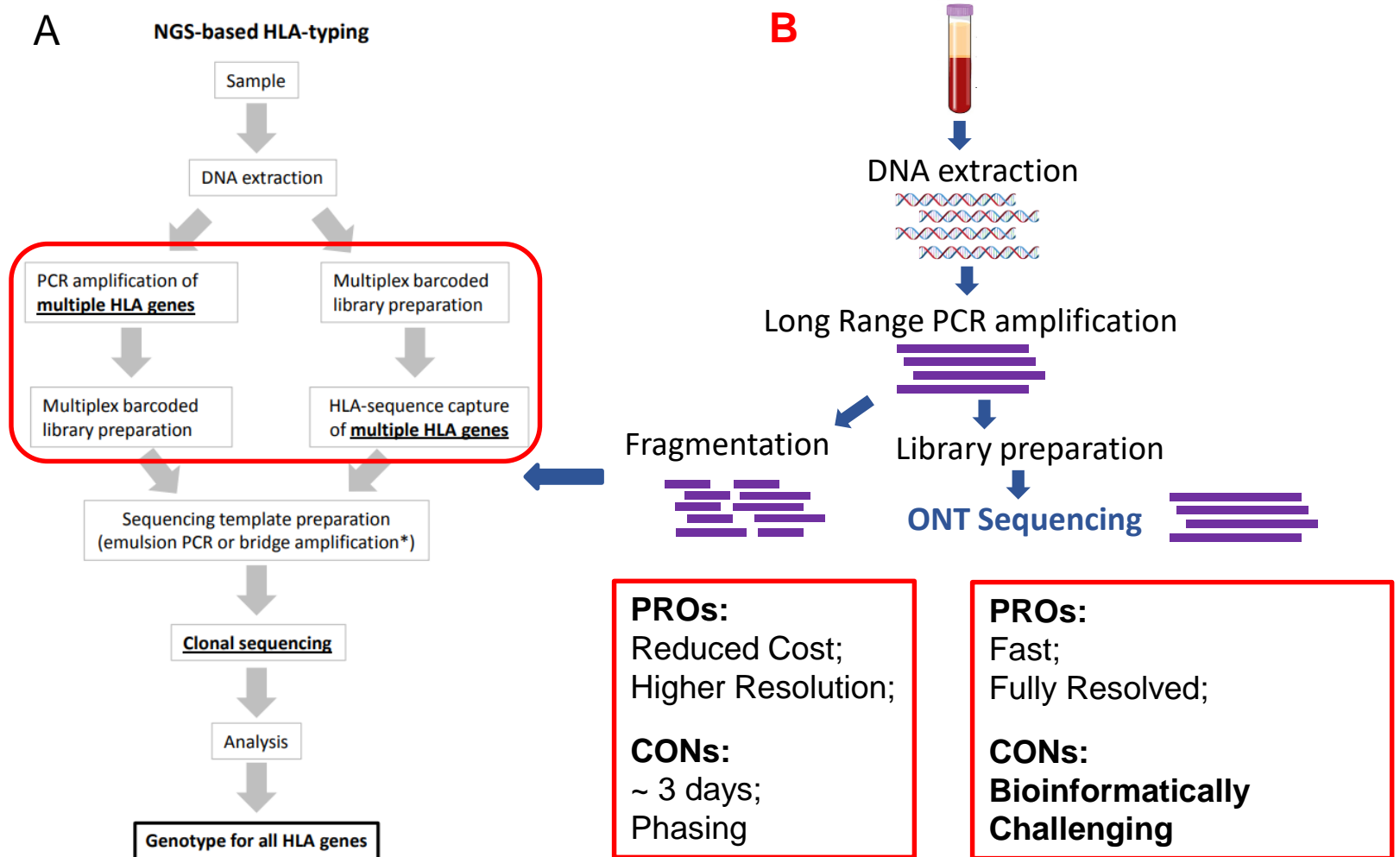
- **C**ytochrome **P**-450 mixed-function oxidase (CYP)
 - abundant in the **liver** and **other organs**
 - responsible for the metabolism of many drugs and environmental chemicals (antiarrhythmics, adrenoceptor antagonists, tricyclic antidepressants)
- Genetic polymorphisms of **CYP2D6**, **CYP2C9**, and **CYP2C19** have been well studied



Examples of Haplotype-Drug Associations

Gene	Haplotypes or genotypes	Phenotypes	Related Drug	Clinical Interpretation
ATM	WT/WT	rs11212617 CC genotype	Metformin	NORMAL RESPONSE
CYP1A2	*1A/*1F	Ultrarapid Metabolizer	Cyclobenzaprine	INCREASE DOSE
CYP2A6	*1/*1	Normal Metabolizer	Nicotine	NORMAL RESPONSE
CYP2C19	*1/*2	Intermediate Metabolizer	Sertraline	USE CAUTION
CYP2C9	*1/*1	Normal Metabolizer	Warfarin	NORMAL DOSE
CYP2D6	*4/*35	Intermediate Metabolizer	Codeine	CONSIDER ALTERNATIVES
CYP3A4	*1A/*1B	Intermediate Metabolizer	Alfentanil	DECREASE DOSE
CYP3A5	*1A/*3A	Expresser	Sirolimus	NORMAL RESPONSE
CYP4F2	*1/*1	Normal Metabolizer	Phenprocoumon	NORMAL RESPONSE
DDRGK1	WT/c.510+364T>G	rs6051639 AC genotype	Ribavirin	USE CAUTION
DPYD	*5/*9A/c.496A>G/IVS10-15T>C	Normal Metabolizer	Capecitabine	NORMAL RESPONSE
F2	WT/WT	Wild Type	Oral-Contraceptives	NORMAL RESPONSE
F5	WT/WT	Non Factor V Leiden Carrier	Eltrombopag	NORMAL RESPONSE
G6PD	WT/WT	Normal G6PD Efficiency	Chlorpropamide	NORMAL RESPONSE
HLA-B	15:02/08:20	15:02	carbamazepine	CONSIDER ALTERNATIVES

Platform for full genotyping of HLA loci



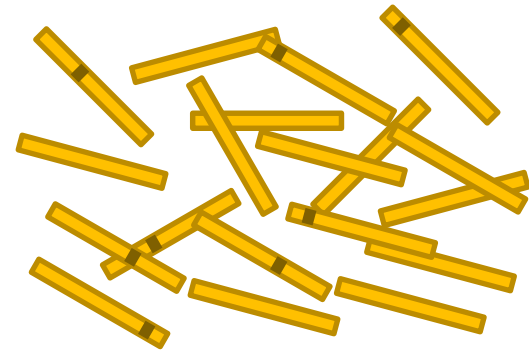
High polymorphic genes associated with drugs efficacy and adverse events



Pharmacogenomic Biomarkers in Drug Labeling					
ALK	CYP2B6	F2	HPRT1	NPM1	ROS1
BCHE	CYP2C19	F5	IFNL3	PDGFRA	SERPINC1
BCR-ABL1	CYP2C9	FIP1L1-PDGFR	IL12A	PDGFRB	SLCO1B1
BRAF	CYP2D6	FLT3	IL12B	PGR	TPMT
BRCA	CYP3A5	G6PD	IL23A	PML-RARA	TPP1
CASR	DMD	GALNS	IL2RA	POLG	UGT1A1
CD274	DPYD	HLA-A	KIT	PROC	VKORC1
CFTR	EGFR	HLA-B	MS4A1	PROS1	
CYB5R	ERBB2	HLA-DQA1	MYCN	RAS	
CYP1A2	ESR	HLA-DRB1	NAGS	RET	

<https://www.fda.gov/Drugs/ScienceResearch/ucm572698.htm>

Long read & Short read



ONT

low (Q10)

Very Long

1kb ~ 882kb¹

Quality


Read Length

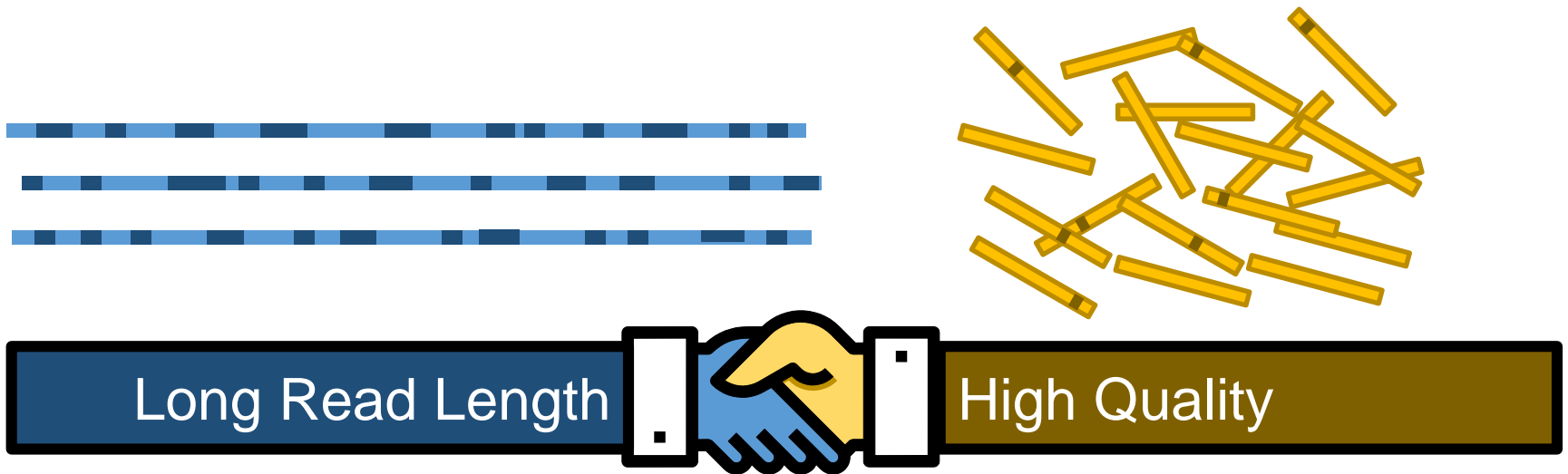
NGS

Very High (Q30)

short

0.5 kb

 The deep color region meant bad quality region.




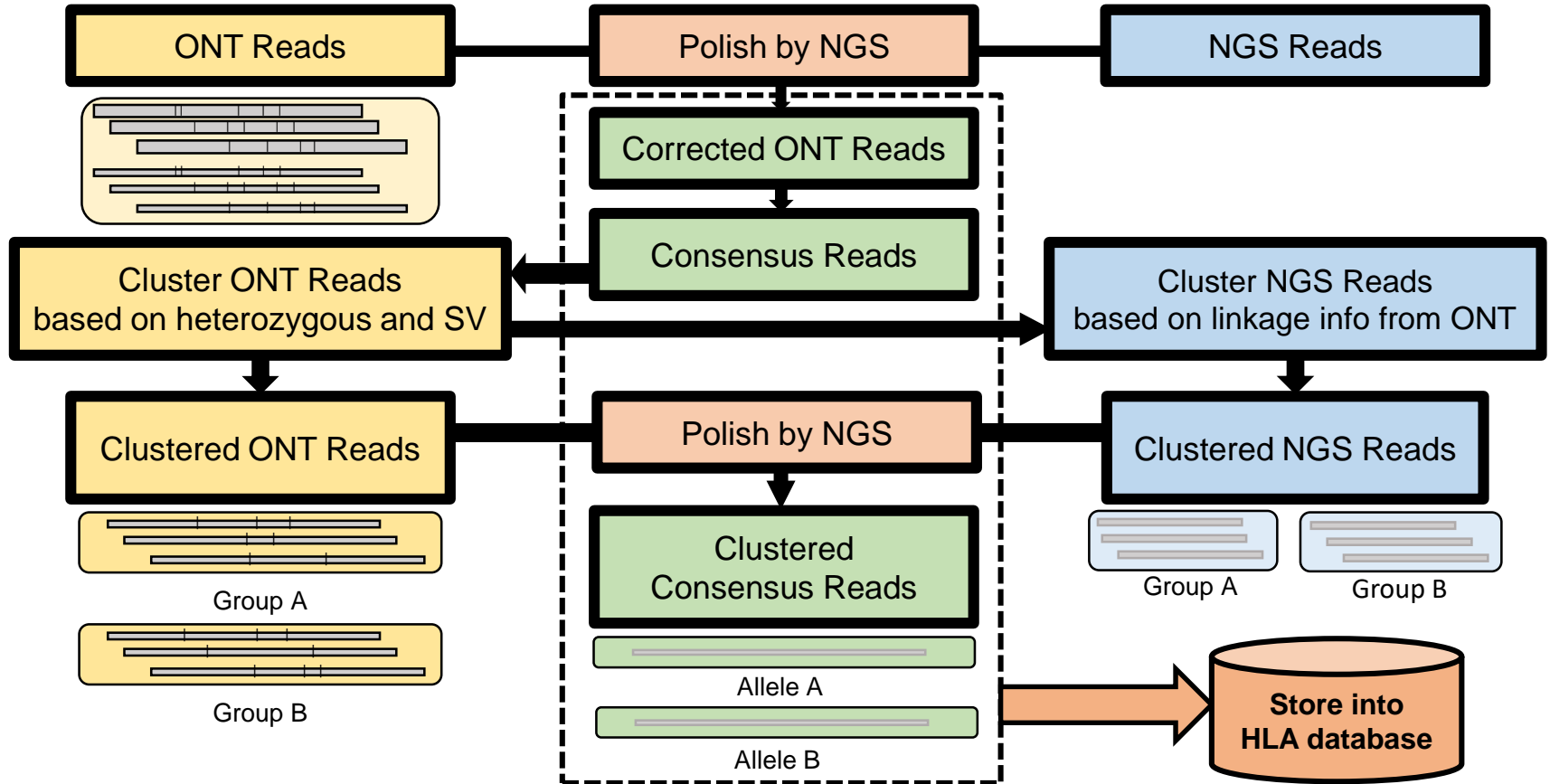
Hybrid correction (Polished)



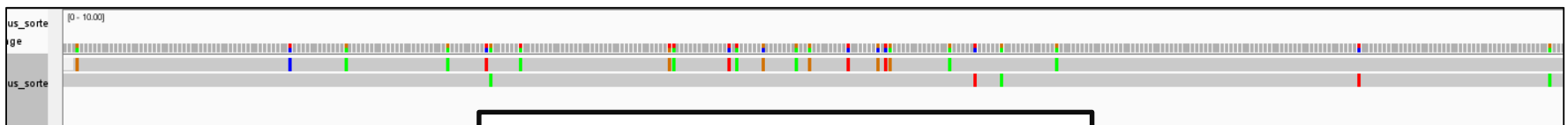
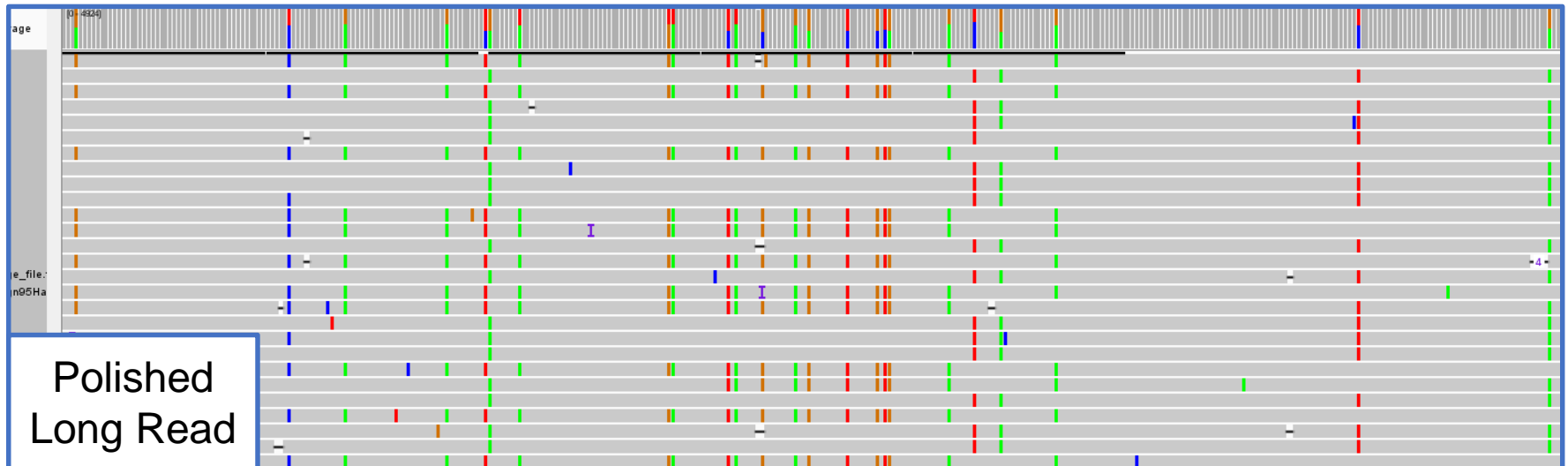
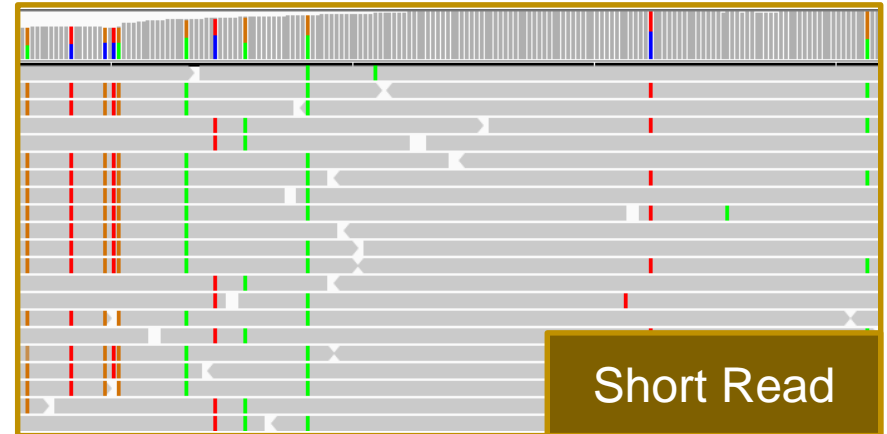
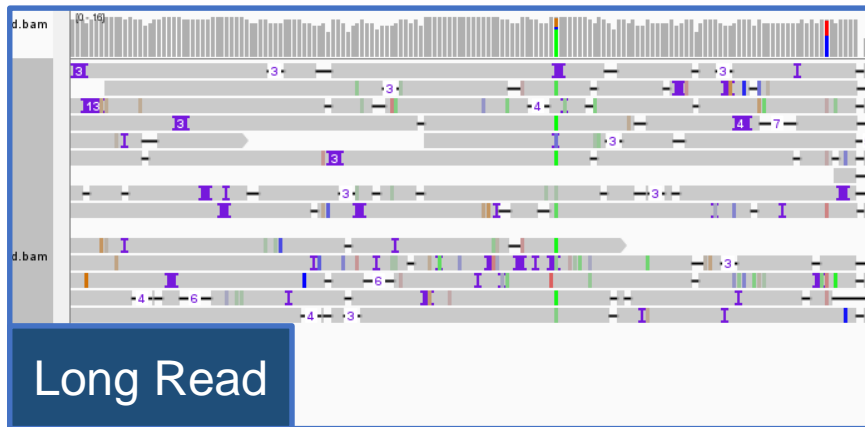
High quality & long read length



Sequencing Analysis Work Flow



Construction of HLA-A alleles of HanB



Top alignment to IMGT: The allele type of HLA-A, HLA-B and HLA-C form HanB

HLA-A*02:10

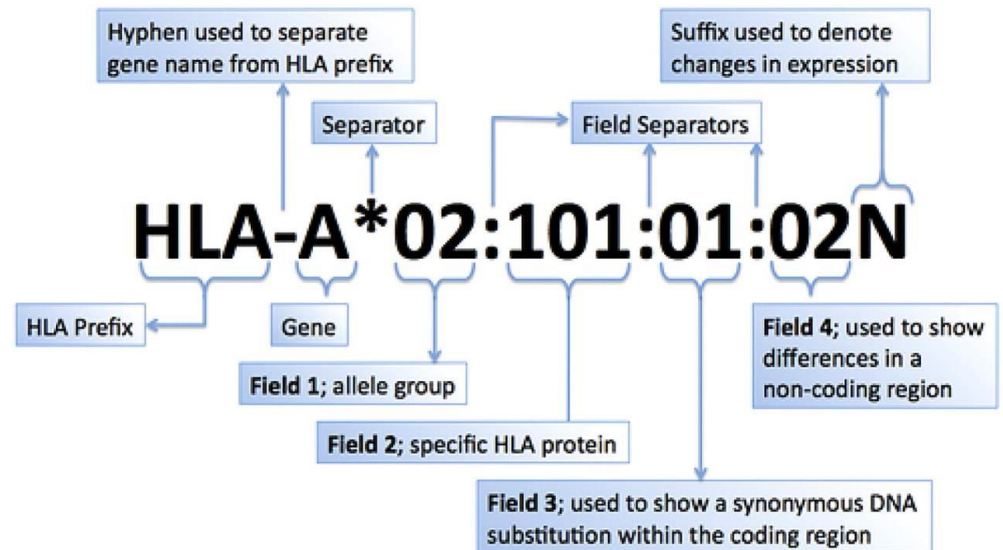
HLA-A*24:02:01:01

HLA-B*15:32:01

HLA-B*40:06:01:01

HLA-C*04:01:01:14

HLA-C*08:01:01:01



© SGE Marsh 04/10

Constructing Database in Different Populations



SCIENTIFIC REPORTS

OPEN Mapping the genetic diversity of HLA haplotypes in the Japanese populations

Received: 17 July 2015
Accepted: 06 November 2015
Published: 09 December 2015

Woei-Yuh Saw^{1,2}, Xuanyao Liu^{1,3}, Ci Tomohiro Katsuya⁴, Ryosuke Kimu Ken Yamamoto^{1,5}, Mitsuhiro Yokota Yik-Ying Teo^{1,2,3,5,13,*} & Norihiro Ka

27:1597–1607 Published by Cold Spring Harbor Laboratory Press; ISSN 1088-9051/17; www.genome.org

Genome Research 1597
www.genome.org

“Several studies^{47–49} have reported the risk of inaccuracies and confounding in genetic association studies in populations even with relatively small genetic differences. In this line, based on our data, we can further advocate caution in using a generic Japanese panel (e.g., JPT in the HapMap) for **imputation of SNPs** and HLA alleles in samples from Okinawa Prefecture.”

Sci Rep. 2015 Dec 9;5:17855.

Assembly and analysis of 100 full MHC haplotypes from the Danish population



Jacob M. Jensen,¹ Palle Villesen,^{1,2} Rune M. Friberg,¹ The Danish Pan-Genome Consortium,⁵ Thomas Mailund,¹ Søren Besenbacher,^{1,3} and Mikkel H. Schierup^{1,4}

¹Bioinformatics Research Centre, Aarhus University, 8000 Aarhus C., Denmark; ²Department of Clinical Medicine, Aarhus University, 8200 Aarhus N., Denmark; ³Department of Molecular Medicine, Aarhus University Hospital, Skejby, 8200 Aarhus N., Denmark; ⁴Department of Bioscience, Aarhus University, 8000 Aarhus C., Denmark

“we reconstruct full MHC haplotypes from de novo assembled trios.... We report 100 full MHC haplotypes and call a large set of structural variants in the regions for future use in **imputation with GWAS data**.”

Genome Res. 2017 Sep;27(9):1597-1607.



Journal of Clinical & Cellular Immunology

Gowda et al., J Clin Cell Immunol 2016, 7:2
<http://dx.doi.org/10.4172/2155-9899.1000399>

Research Article

Open Access

Comparative Analyses of Low, Medium and High-Resolution HLA Typing Technologies for Human Populations

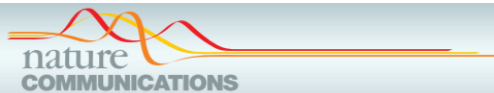
Malali Gowda^{1,2*}, Sheetal Ambaradar¹, Nutan Dighe³, Ashwini Manjunath¹, Chandana Shankaralingu¹, Pradeep Hirannaiah¹, John Harting⁴, Swati Ranade⁴, Latha Jagannathan³ and Sudhir Krishna²

¹Next Generation Genomics Laboratory, Centre for Cellular and Molecular Platform, National Centre for Biological Sciences, TIFR Bangalore, India

J Clin Cell Immunol 2016, 7:2

“This is the first case study of HLA typing using second and third generation NGS technologies for an Indian population. The PacBio platform is a promising platform for large-scale HLA typing for **establishing an HLA database for the untapped ethnic populations of India**.”

Significant Implications of HLA Haplotyping



MHC matching improves engraftment of iPSC-derived neurons in non-human primates

Asuka Morizane¹, Tetsuhiro Kikuchi¹, Takuya Hayashi², Hiroshi Mizuma², Sayuki Takara², Hisashi Doi², Aya Mawatari², Matthew F. Glasser³, Takashi Shiina⁴, Hirohito Ishigaki⁵, Yasushi Itoh⁵, Keisuke Okita⁶, Emi Yamasaki¹, Daisuke Doi¹, Hirohito Onoe^{2,7}, Kazumasa Ogasawara⁵, Shinya Yamanaka^{6,8} & Jun Takahashi^{1,9}

Nat Commun. 2017 Aug 30;8(1):385.

Common Allele/Haplotyping

BIOINFORMATION

Discovery at the interface of physical and biological sciences

Bioinformatics. 2017; 13(3): 94–100.
Published online 2017 Mar 31. doi: [10.6026/97320630013094](https://doi.org/10.6026/97320630013094)

PMCID: PMC5450251

T-cell epitopes predicted from the Nucleocapsid protein of Sin Nombre virus restricted to 30 HLA alleles common to the North American population

Sathish Sankar,^{1,*} Mageshbabu Ramamurthy,¹ Balaji Nandagopal,¹ and Gopalan Sridharan¹

Bioinformatics. 2017; 13(3): 94–100.

REPORTS

Cite as: D. Chowell *et al.*, *Science* 10.1126/science.aao4572 (2017).

Allelic Expression

Unique Allelic eQTL Clusters in Human MHC Haplotypes



Tze Hau Lam,^{*} Meixin Shen,^{*} Matthew Zirui Tay,[†] and Ee Chee Ren^{*,†,1}

^{*}Singapore Immunology Network, A*STAR, Singapore 138648, [†]Department of Molecular Genetics, Duke University, Durham, North Carolina 27710, and [‡]Department of Microbiology and Immunology, School of Medicine, National University of Singapore, Singapore 117597

Science

Immunotherapy (Antibody based and Cell based)

Patient HLA class I genotype influences cancer response to checkpoint blockade immunotherapy

Diego Chowell,^{1,2} Luc G.T. Morris,^{2,*} Claud M. Grigg,^{**} Jeffrey K. Weber,⁵ Robert M. Samstein,^{1,2} Vladimir Makarov,^{1,2} Fengshen Kuo,^{1,2} Sviatoslav M. Kendall,^{1,2} David Requena,⁶ Nadeem Riaz,^{1,2,7} Benjamin Greenbaum,⁸ James Carroll,⁹ Edward Garon,⁹ David M. Hyman,^{10,14} Ahmet Zehir,¹¹ David Solit,^{1,10,12} Michael Berger,^{1,11,12} Ruhong Zhou,^{5,13} Naiyer A. Rizvi,^{4†} Timothy A. Chan^{1,2,7,14†}

[†]Human Oncology and Pathogenesis Program, Memorial Sloan Kettering Cancer Center, New York, NY 10065, USA. [‡]Immunogenomics and Precision Oncology Platform,

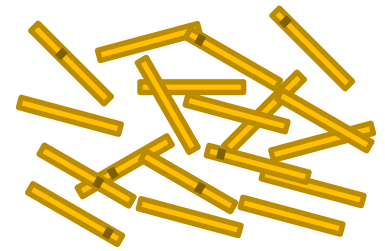
Science. 2018 Feb 2;359(6375):582-587.



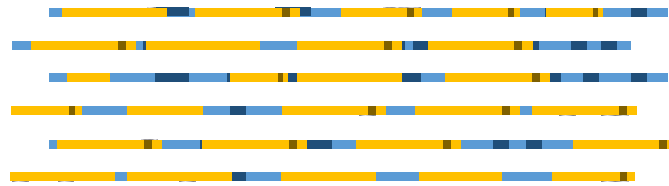
Long Read



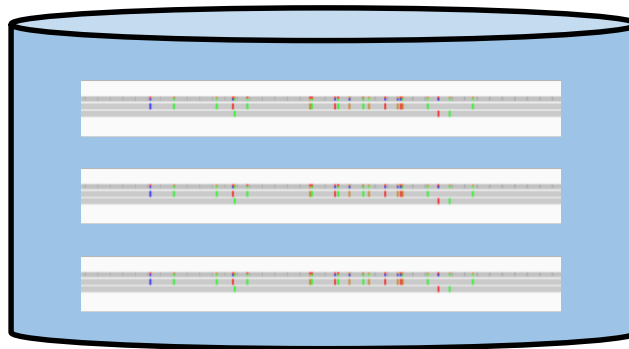
Polished Long Read



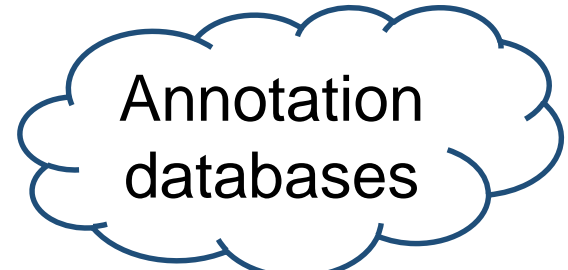
Short Read



Full-length alleles



Alleles database



Haplotype genotyping

Allele genotyping

Precision Medicine

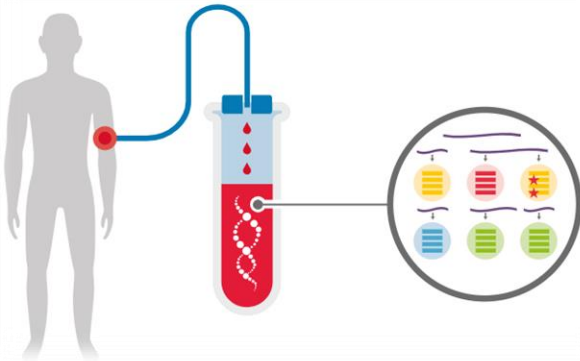
Sequencing Platform Comparison

Platform Comparison	Illumina	Thermo Fisher	Pacific Biosciences	Oxford Nanopore
Sequencing by Synthesis	Yes	Yes	Yes	No
DNA Size Selection	Yes	Yes	Yes	No
Post-library Amplification	Yes	Yes	No, Single Molecule	No, Single Molecule
Detection	Fluorescent Imaging	Ion Semiconductor	Fluorescent Imaging	Ionic Current Change
Sequencing Rate (s/base)	2 – 20 sec	30 sec	0.25 sec	0.002 sec
Running Time	Fixed	Fixed	Fixed	Run & Stop
DNA Sequencing	Yes	Yes	Yes	Yes
Direct DNA Modification Detection	No	No	Yes	Yes
Direct RNA Sequencing	No	No	No	Yes
Read Length	Short, up to 300 bp X 2	Short, up to 600 bp	Long, Average 6-8 Kb	Long, Average 6-30 Kb
Total Reads (M)	4 – 800 (PE)	2 – 130	0.3 – 0.5	0.3 – 0.5
Total Base (Gb)	1.2 - 120	0.3 – 25 /Chip	5 – 8 /SMRT Cell	2 – 10 /Flow Cell
Instrument Cost (USD)	20K – 275K	200 – 300K	350K	1K – 125K

Applications Requiring Throughput (coverage)

- Liquid biopsy for diagnosis and prognosis
- Liquid biopsy for reproductive and genetic health
- Human genome resequencing
- Signal counting applications, such as gene expression profiling
- Detection of DNA mutation (diagnostic testing)

Liquid biopsy



Resequencing



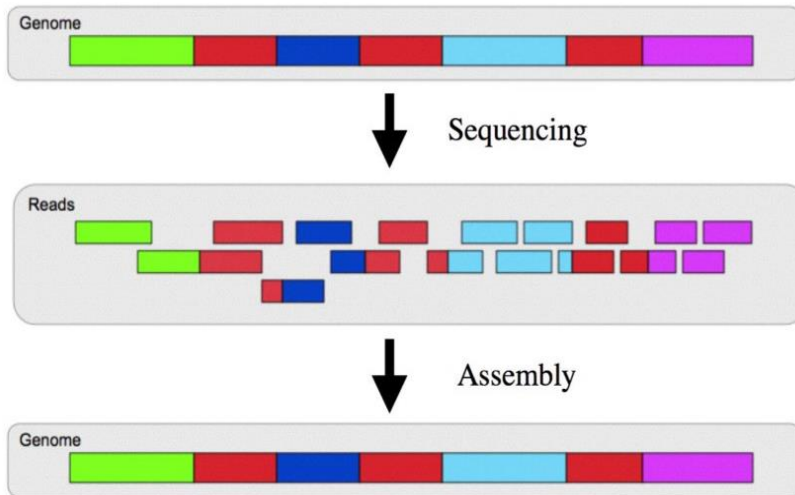
Coverage



Applications Requiring Long Reads

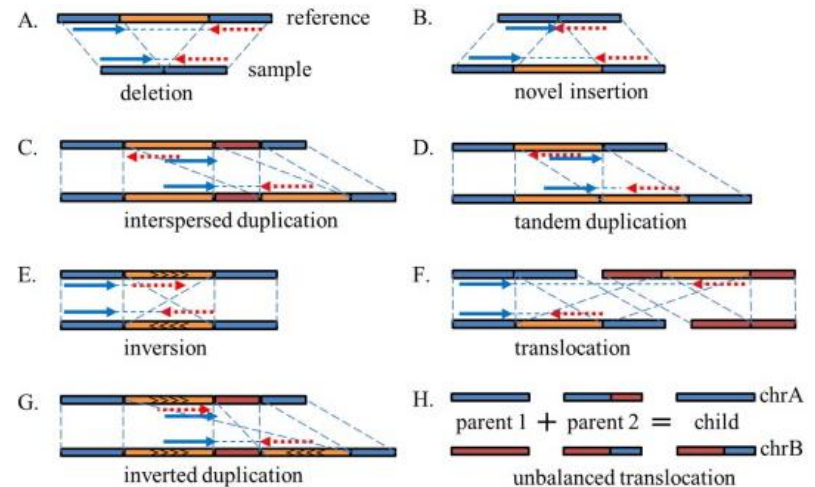
- Genome assembly (New resolution)
- Detection of structural variations
- Full genotyping of highly polymorphic loci

Genome assembly



<https://www.youtube.com/watch?v=5wvGapmA5zM>

Structural variation detection



Methods. 2016 Jun 1;102:36-49.



Thank you for your attention

