

Diagnostic applications of Clonal Sequencing (NGS)

Prof. Graham Taylor PhD FRCPath
Leeds Teaching Hospitals &
Leeds Institute of Molecular Medicine

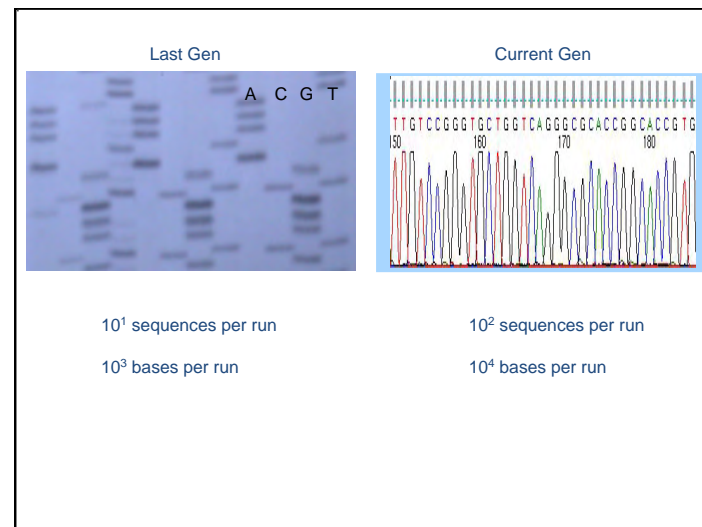
Diagnostic applications of NGS

- Background
- Workflow
- Diagnostic Applications
 - Resequencing
 - PCR
 - Capture
 - CNV-seq
- Summary & Next Steps

Sequencing Goes Large

Roche 454	Illumina Solexa	ABI SOLiD
--------------	--------------------	--------------

Helicos	Pacific Bioscience	Oxford Nanopore
---------	-----------------------	--------------------

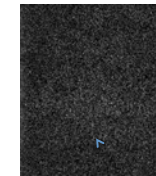
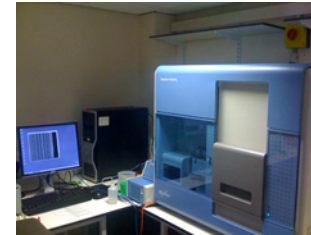


Lessons from the 8q24 study

- Short reads (32 base) can cover long genomic regions because repeats are imperfect and can be aligned
- Long PCR is an effective way to produce template for sequencing
- Coverage is uneven (why?)
- Point mutations are detectable and consistent with conventional sequencing results
- Indels are missing or under-represented
- Cost per base about **50 fold less** with Illumina than Roche

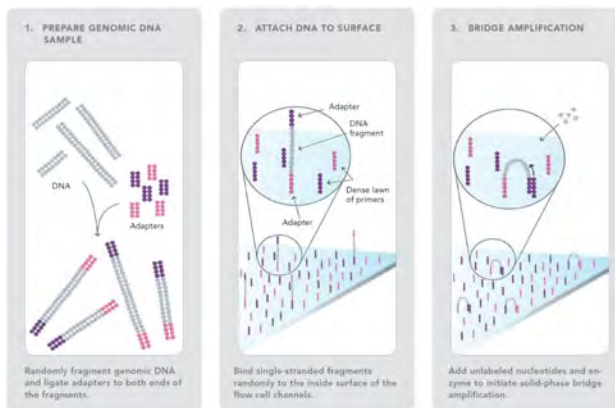
Application of Illumina GAI

Cost: 0.002p per base
 Capacity: >7.0 Gigabase per run
 Simplicity:
 library>cluster
 station>sequence>data

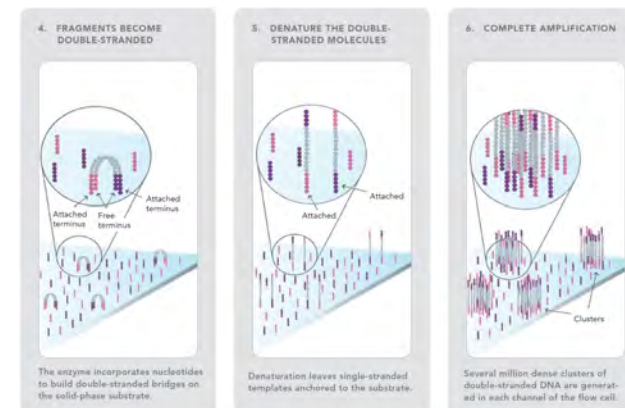


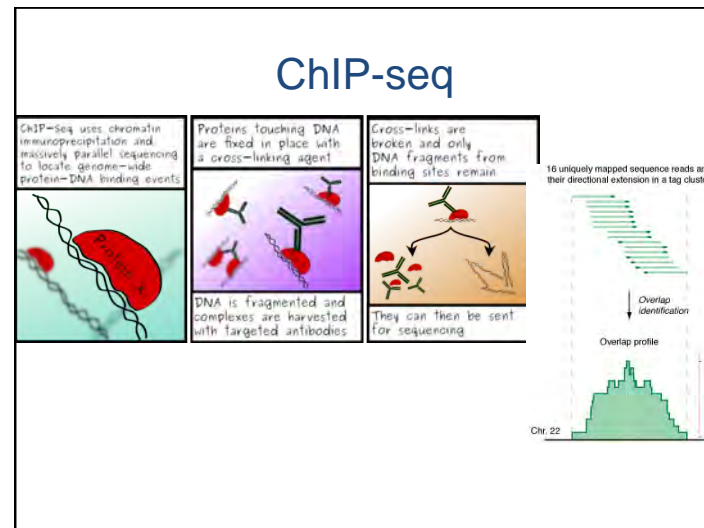
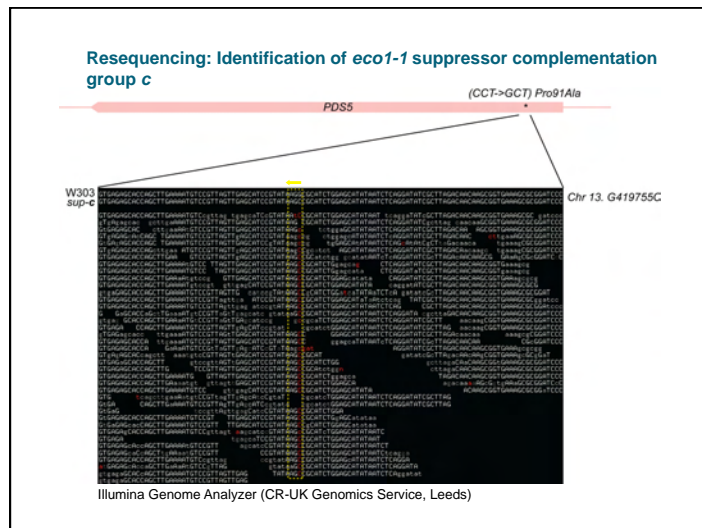
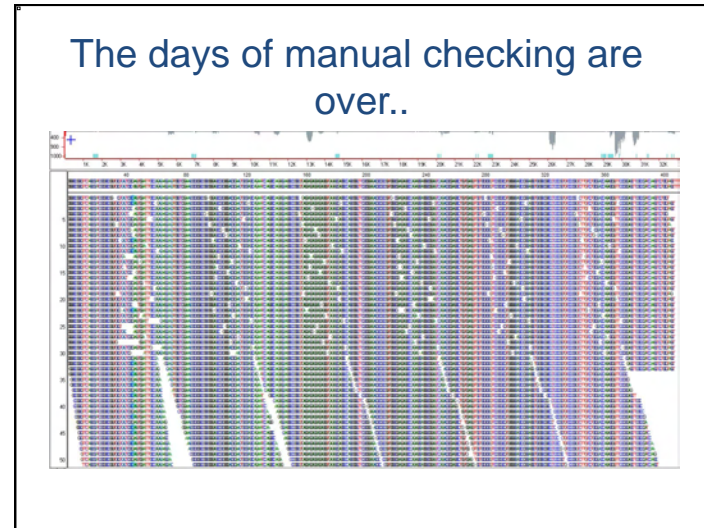
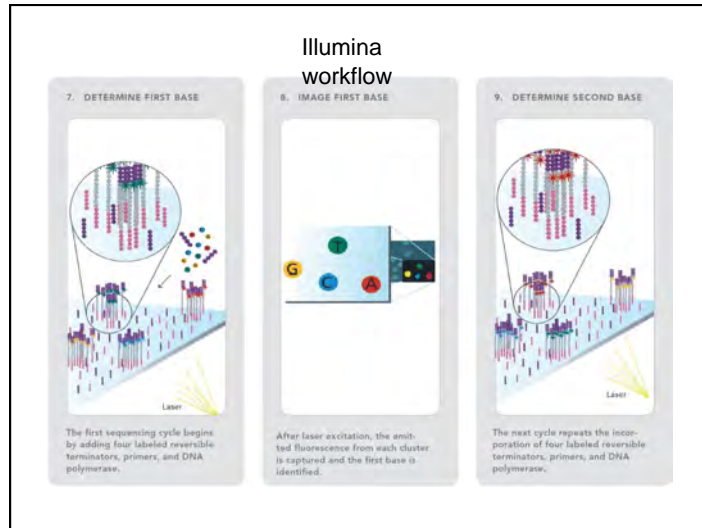
100,000 spots
 (1/100th of one
 flow-cell channel)

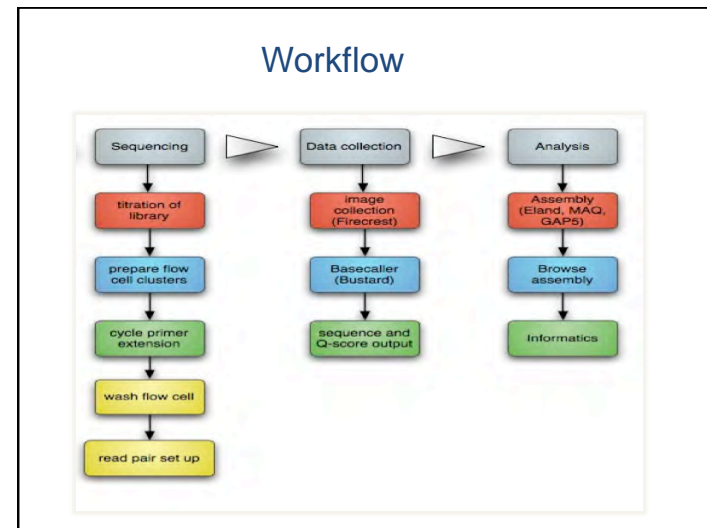
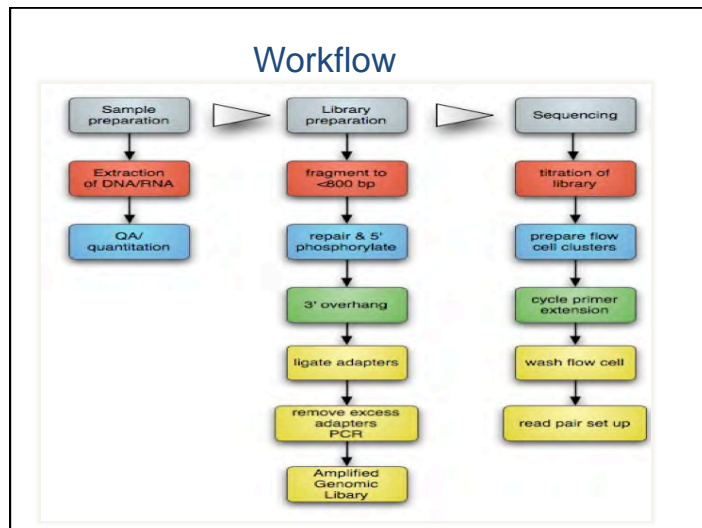
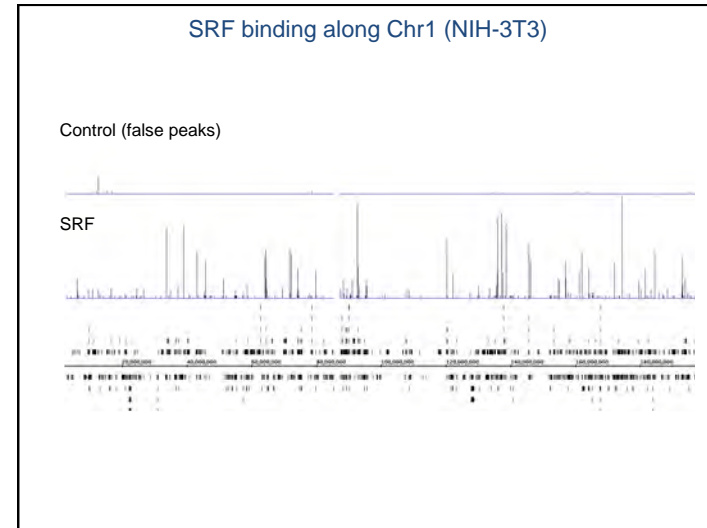
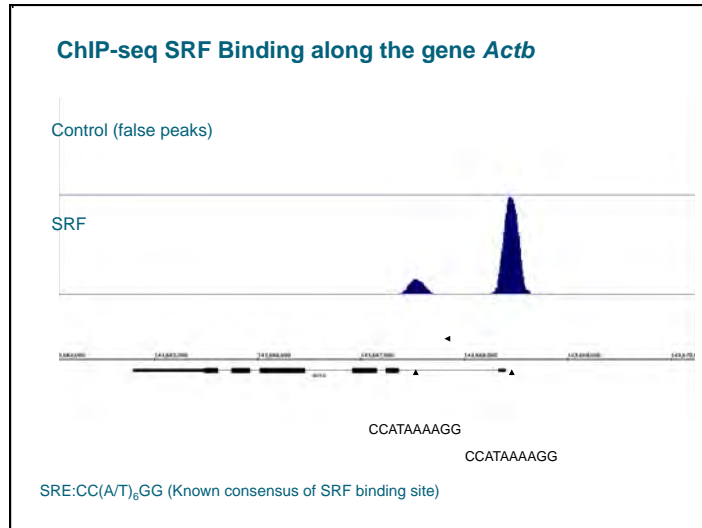
Illumina workflow



Illumina workflow







Software

- Pipeline: Signal>Base calls
 - Illumina: GOAT, Bustard and Firecrest (Unix)
- Assembly, Mapping & Browsing
 - ELAND (Illumina, Unix)
 - MAQ & Maqview (Sanger, Unix)
 - NextGene (Softgenetics, Windows)

What does this mean for clinical diagnostics?

The currency of genetic analysis will become DNA sequence

- Sequence count
- Sequence variation
- Sequence arrangement
- The currency is dropping to commodity prices
- Skill set required is less lab, more informatics biased

How will we use it?

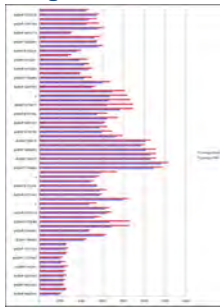
- Test design when DNA sequence costs are trivial
 - Urgent (same day) or non urgent (2 week)
 - Pre and post test costs
 - Syndrome/Gene list driven
 - Hereditary cancers, “Cardiome” “Retinome”
 - “Kinome” : 20 – 200 genes 250-3,000 megabases

Clinical Applications

- Incremental rather than radical
- Mutation Detection in *TP53*, *BRCA1* & *BRCA2*
- Gene capture in autozygosity projects
- Gene dosage & chromosome counting
 - Fixed tissue
- What does it mean for the clinical service?

Coverage and costs

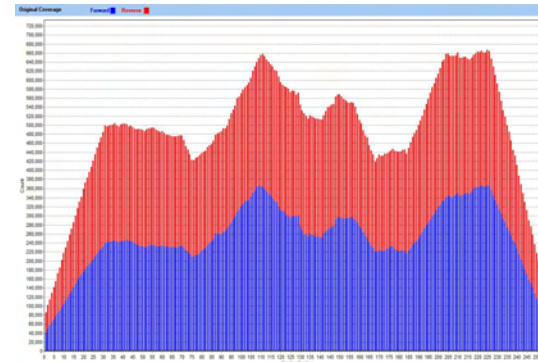
Coverage



Costs

Capital Outlay	220000
Pipeline costs	30000
Total Project Outlay	250000
Annual Running Costs:	
Pay	
1 FTE - SSO	37500
1 FTE - SC2	30000
0.5 FTE Bioinformatics	18000
Total Pay Costs	85500
Non Pay	
Reagent costs	100000
Maintenance costs	20000
Depreciation	44000
Total Non Pay Costs	164000
Total Annual Running Costs	249500

Coverage can be wide (e.g.exome) or deep (e.g. pooled samples)



Sequence wide or deep to achieve cost efficiency

- Sample pooling
 - Ingman & Gyllensten EJHG 17 383-6 2009
- Sample tagging
 - Craig *et al* Nat Methods 5, 887-93, 2008
 - 3 base
 - 6 base tag = 4^5 (1024) variants
 - Illumina kit

The case for gene-centric analysis

Coverage

- 1 flow cell = 7 channels
- 1 channel = 10 million reads
- 1 read = 36 – 200 bases
- 1 channel = 0.3 – 2 Gigabases
- 1 flow cell = 4 haploid genome equivalents
- At mean coverage of 100, one channel can cover 5-20 Megabases

Costs

- 1 channel costs between £600-£1,200 in reagents
- 1 base costs between 2/10,000 pence – 5/100,000 pence
- At mean coverage of 100, one base costs 1/50 – 5/10,000 pence

Target selection methods

Sequence everything and select later in software

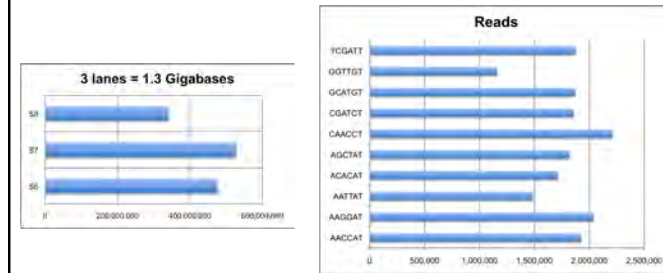
OR

Long PCR

Gene Collector

Genome Partitioning

6 base tags 51 base reads



Long PCR

- 8q24
- TP53
- CDKN2A
- BRCA1 and 2 exon 11

Current software releases not ideal

- ELAND/CASAVA/Genome Studio: supplied with the machine.
- Academic
- Commercial
 - NextGene
 - CLC Workbench
 - DNASTar
- Illuminator
 - <http://dna.leeds.ac.uk/illuminator/>

Additional software used

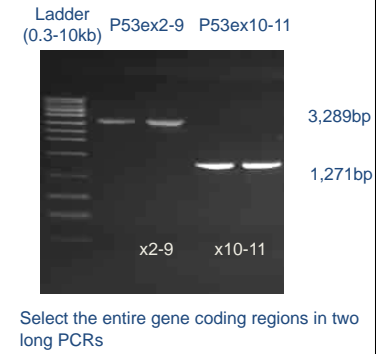
<http://www.softgenetics.com/>

<http://dna.leeds.ac.uk/sequencing/>

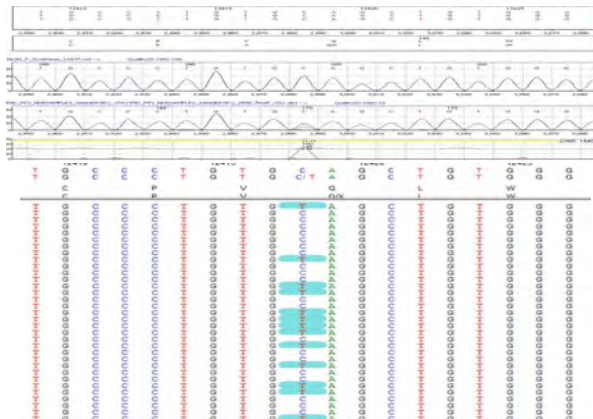


Proof of principle

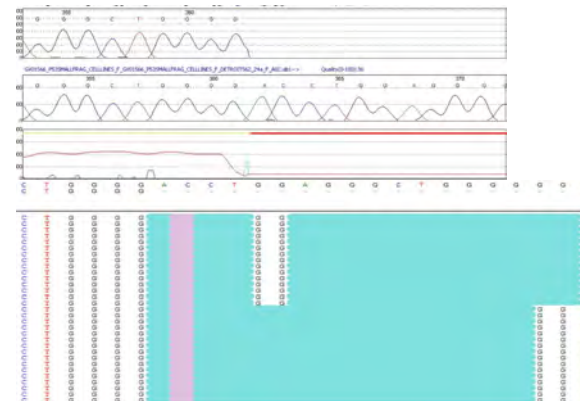
- *TP53*: mutated in 50% somatic cancers
- Small gene (4,510 bases of *TP53*, including 1,172 bases of coding sequence)
- Cell lines and cases available



Point mutation (p.144Q>QX)



16 base insertion-deletion

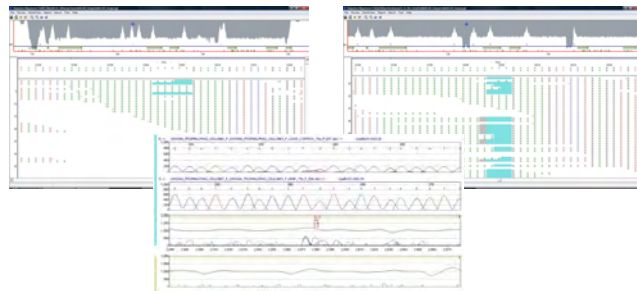


False +ves: slippage at (A)₁₈ Poor quality Sanger sequencing

Manageable number of false positives, very few in or near exons

30 base reads

45 base reads



TP53 Results

4 cell lines in duplicate

- 5 coding variants reported
- Duplicate runs fully concordant
- No false negatives
- All pathogenic mutations identified in COSMIC or by Sanger sequencing were detected

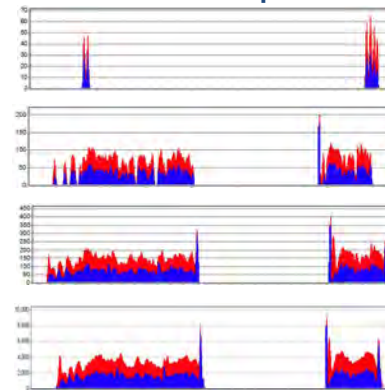
10 cases in duplicate

- 45 variants reported within 50 bases of exons 2-11
- Duplicate runs fully concordant
- No false negatives
- All pathogenic mutations previously seen by Sanger sequencing were identified

TP53 Coding Variants

Case	Amino Acid Change	Variant Call	Reported Mutation
1	72P>R	c.[215C>G]+[215C>G]	
	144Q>QX	c.[410C>T]+[+]	c.A30C>T (p.Gln144X)
2	377T>TP	c.[1129A>C]+[+]	
	72P>R	c.[215C>G]+[215C>G]	
3	248R>QR	c.[743G>A]+[+]	c.743G>A (p.Arg248Gln)
	72P>R	c.[215C>G]+[215C>G]	
4	248R>QR	c.[743G>A]+[+]	c.743G>A (p.Arg248Gln)
	371S>SS	c.[1113C>C]+[1113C>C]	
5	377T>TP	c.[1129A>C]+[+]	
	72P>R	c.[215C>G]+[215C>G]	
6	158R>HR	c.[473G>A]+[+]	c.473G>A (p.Arg158His)
	377T>TP	c.[1129A>C]+[+]	
7	72P>PR	c.[215C>G]+[+]	
	158R>PR	c.[473G>C]+[+]	c.473G>C (p.Arg158Pro)
8	72P>R	c.[215C>G]+[215C>G]	
	256T>TA	c.[766A>G]+[+]	c.766A>G (p.Trp256Asp)
9	377T>TP	c.[1129A>C]+[+]	
	72P>PR	c.[215C>G]+[+]	
10	273R>HR	c.[818G>A]+[+]	c.818G>A (p.Arg273His)
	377T>TP	c.[1129A>C]+[+]	
11	72P>R	c.[215C>G]+[215C>G]	
	242C>XC	c.[776C>A]+[+]	c.726C>A (p.Cys242X)
12	377T>TP	c.[1129A>C]+[+]	
	125T>TR	c.[1374C>G]+[+]	c.374C>G (p.Trp125Arg)
13	377T>TP	c.[1129A>C]+[+]	
	72P>PR	c.[215C>G]+[+]	
14	158R>HR	c.[473G>A]+[+]	c.473G>A (p.Arg158His)

Read Depth vs. Detection



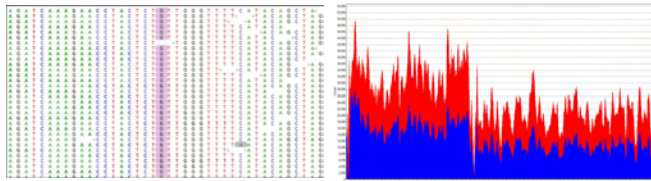
Reads	Samples	Detection
5,000	1,000	0/11
10,000	500	7/11
20,000	250	11/11
500,000	10	11/11

Mutation Call	Amino Acid Change
c.[74+38C>G]+[74+38C>G]	
c.96+41delACTAGGCTGGGGGG	
c.[215C>G]+[215C>G]	72P>R
c.[375-283T>C]+[375-283T>C]	
c.375-160delAAA	
-	
-	
c.[375-91G>A]+[375-91G>A]	
c.[451C>T]+[451C>T]	151P>S
c.[672+62A>G]+[672+62A>G]	
c.[672-66T>C]+[+]	
c.[672-36G>C]+[672-36G>C]	

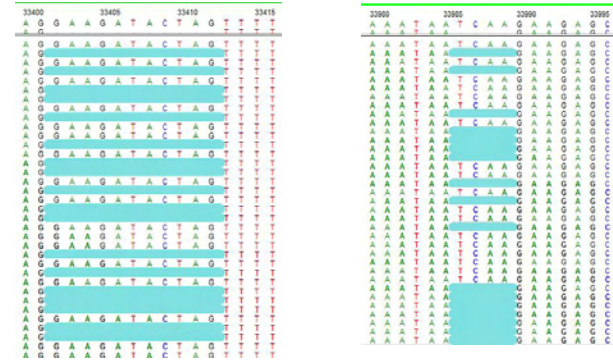
BRCA1 and *BRCA2* exon 11 control

Typical sequence with SNP

Average depth: 25,000



BRCA1 exon 11



BRCA1 exon 11 and *BRCA2* exon 11 summary

- 13 of 13 mutations plus known SNPs were detected and reported automatically, even though present as 1/8 or 1/6 of a sample mixture
- Even at this crude pilot level, the cost for complete sequence screening is competitive with existing techniques, and faster

BRCA1 and *BRCA2*

- 65,159 bases from the *BRCA1* and *BRCA2* genes in 22 amplicons
- 16,032 coding bases
- 55 previously sequenced cases
- Average sequencing coverage 592
- One run repeated 6 times
- One control sample included in all runs

Gene collector

Multiplex amplification of all coding sequences within 10 cancer genes by Gene-Collector
 Simon Friedlacker*, Jolien Bandy, Fredrik Dahl, Angela Chu, Harbir Ji, Kaitoko Wachi and Donald W. Brown
 Swedish Cancer Technology Center, SBC-C, 218 Gårdsgränd, Stockholm, S-141 86, Sweden
 *Present address: 218 Gårdsgränd, S-141 86, Stockholm, Sweden

- Multiplex PCR (pfu polymerase), blunt-ended products suitable for ligation by circularization
- Collector probes guide circularization, closed circles formed by thermostable ligase
- Enrichment of circular DNA by exonuclease treatment and rolling circle amplification
- Risk of artefacts from RCA

Gene collector

Addresses the analysis of short fragments by Illumina

- XPC
- KRAS/BRAF in fixed tissue
- BRCA 1 and 2: towards complete analysis of BRCA sequence by long PCR and gene collector

Mutation Detection

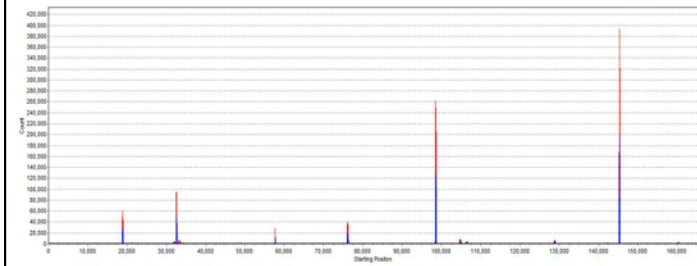
Reference Position	Segment Description	Segment Position	Reference Nucleotide	Coverage	A (%)	C (%)	G (%)	T (%)	Ins (%)	Del (%)
50	xxPC exon1	50	C	28	0.00	0.00	0.00	0.00	78.57	0.00
94	xxPC exon1	94	G	50	0.00	56.90	-56.90	0.00	0.00	0.00
1150	xxPC exon 5	1184	T	45	0.00	0.00	0.00	0.00	100.00	0.00
1863	xxPC exon 9_10_11	379	G	44	0.00	9.09	-100.00	0.00	90.91	0.00
1935	xxPC exon 9_10_11	451	G	20	0.00	100.00	-100.00	0.00	0.00	0.00
1936	xxPC exon 9_10_11	452	A	20	-100.00	0.00	0.00	0.00	100.00	0.00
2017	xxPC exon 9_10_11	523	C	111	0.00	-45.05	0.00	45.05	0.00	0.00
2452	xxPC exon 12_13	34	T	35	100.00	0.00	0.00	-100.00	0.00	0.00
2564	xxPC exon 12_13	148	C	42	0.00	-100.00	0.00	0.00	100.00	0.00

Ras Raf Mutation Report (c.38G>A)

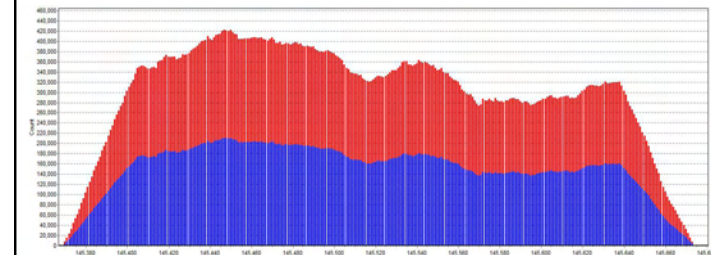
Control

Tumour

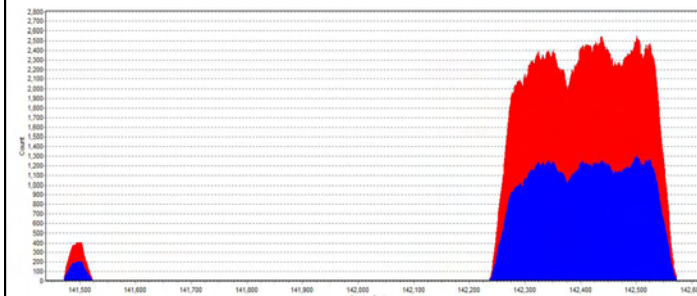
BRCA gene collector



BRCA gene collector



BRCA gene collector



BRCA Mutation Report

Segment Description	Position	Ref Base	Coverage	A(%)	C(%)	G(%)	T(%)	a a change
BRCA1-NCBI-NP_009225.2.gbk	18858	G	55240	0	0	98.2	1.8	
BRCA1-NCBI-NP_009225.2.gbk	19077	T	44629	0	2.22	0	97.8	
BRCA1-NCBI-NP_009225.2.gbk	19118	G	41437	0	0	99.4	0.65	
BRCA1-NCBI-NP_009225.2.gbk	33789	A	3043	91.3	0	8.68	0	1290K>KE
BRCA1-NCBI-NP_009225.2.gbk	33806	G	3341	21.1	0	78.9	0	1295L>LL
BRCA1-NCBI-NP_009225.2.gbk	33913	A	2679	84.9	0	0	15.1	1331Q>QL
BRCA1-NCBI-NP_009225.2.gbk	57763	G	29505	2.25	0	97.8	0	1665V>MV

Matched vs. Unmatched reads

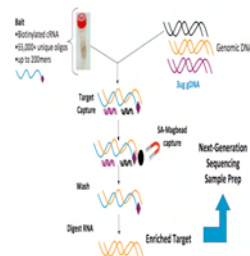
Sample	Matched	Unmatched
TP53 long PCR	20,763	465 (2%)
TP53 long PCR + templphi	17,097	31,646 (65%)
XPC gene collector	15,411	49,860 (76%)
Kras/Braf (FFPE) gene collector	6,368	26,669 (81%)
BRCA1 & BRCA2	21,276	12,432 (37%)

Summary of Gene Collector

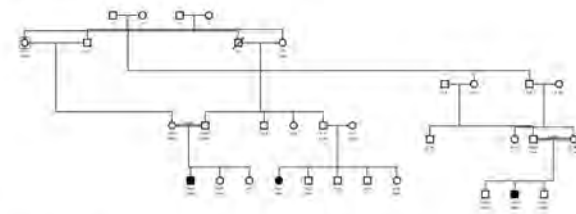
- Converts short amplicons to Illumina-readable format
- Identifies mutations
- Some loss of useful reads compared with long PCR
- Difficult to achieve even amplification

Genome partitioning

- No PCR
- Genome Fragmentation
- Linker and Tag Ligation
- Oligo Capture
- Sequence



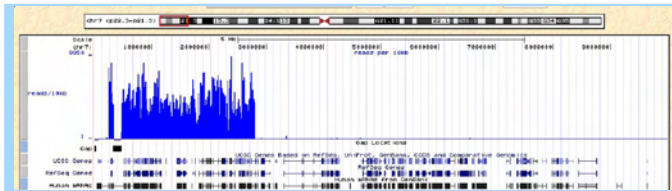
Primary Ciliary Dyskinesia



Defect of the cilia affecting lungs, sinuses, ears and nose.
Can affect organ position- situs inversus.

Autosomal recessive inheritance.
SNP analysis revealed homozygous target region on tip of chr7p.
Contains 32 genes with 300 exons.
Designed custom Agilent capture chip for 2.7Mb region.

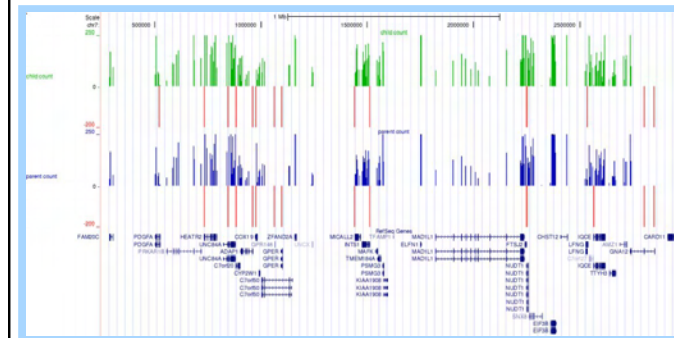
Chr7p target region enrichment



17M reads aligned to genome, 2.2M aligned to target region

Target region is ~3Mb or ~0.1% of genome, giving enrichment of 130-fold

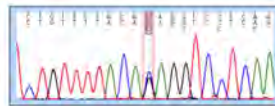
Read depth across the exons



Identification of a splice site mutation

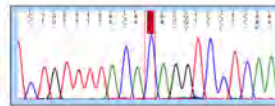
Parent (heterozygous carrier)

Boundary change at 786115 seen in 3 reads of 37
 Boundary change at 786307 seen in 2 reads of 54
 Boundary change at 791679 seen in 7 reads of 15
 Boundary change at 838687 seen in 1 reads of 63
 Boundary change at 838764 seen in 2 reads of 62



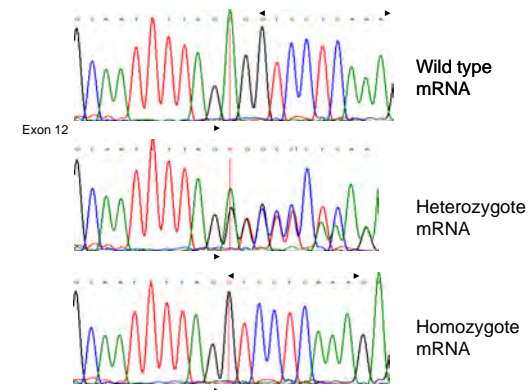
Child (homozygous affected)

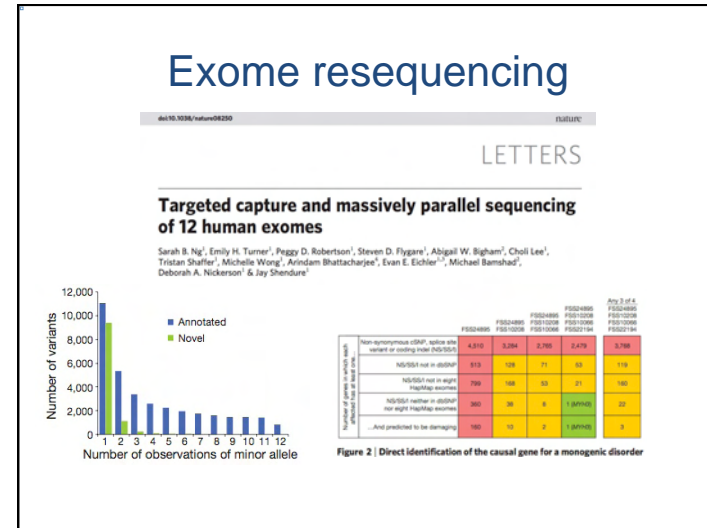
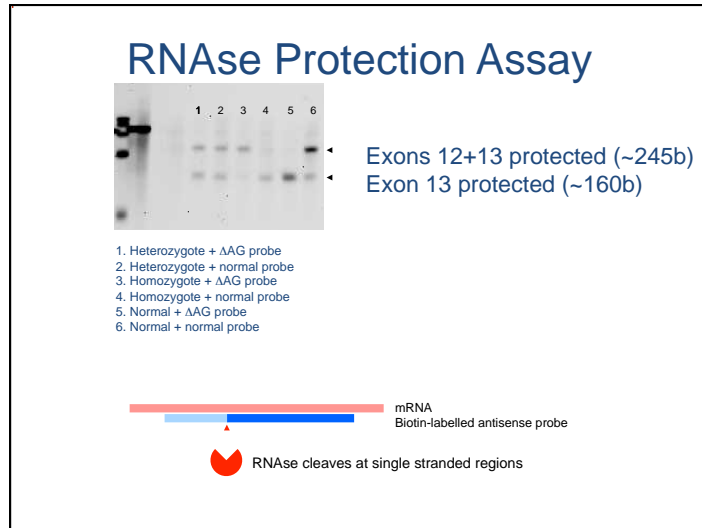
Boundary change at 786115 seen in 2 reads of 68
 Boundary change at 786307 seen in 1 reads of 96
 Boundary change at 791679 seen in 36 reads of 36
 Boundary change at 838687 seen in 2 reads of 116
 Boundary change at 838764 seen in 5 reads of 170



Mutation segregates with the condition
 Affects the splicing of the mRNA
 Is not present in 200 ethnically matched controls
 Is the only pathogenic change identified in the target region
 Highly conserved

Splice site mutation activates cryptic site resulting in 2b deletion in mRNA





Exome Sequencing

- “Targetted capture and massively parallel sequencing of 12 human exomes” Ng *et al.* Nature Sept 2009
 - Captured the entire human exome using Agilent arrays (1% of genome)
 - Generated 30Mb sequence per individual spread over 180000 exons
 - Sequenced 8 HapMap individuals
 - Sequenced 4 unrelated individuals with Freeman-Sheldon syndrome, an autosomal dominant disorder due to mutations in MYH3.
 - Averaged 17272 cSNPs per individual

Exome data

	1	2	3	4
Number of shared variants	4510	3284	2765	2479
Variants not present in dbSNP	513	128	71	53
Variants not present in 8 exomes	799	168	53	21
Variants not present in 8 exomes and dbSNP	360	38	8	1 (MYH3)

But: 4% of exons (7000) were not sequenced deeply enough for analysis

Sub-Genomic Targetting

- Focus on regions of interest
- Higher coverage
- Easier informatics
- Fewer “unwanted” genotypes
- Target on the basis of phenotype or pathways
- Compare with exome targetting
- *Work in progress...*

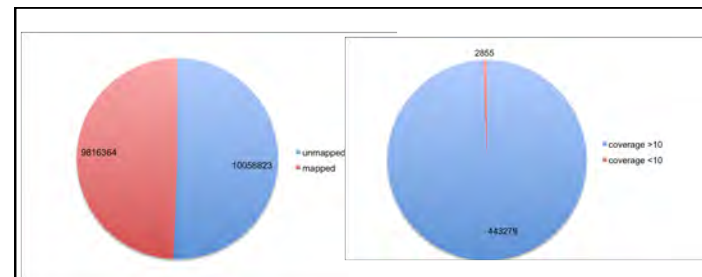
Cardiopathy

- 55 genes
- 600Kb coding region
- Simple repeats were retained in the baits

Class of repeat	Number of each repeat type identified	Percent (%)	Average size (bp)
SINE	015	013.6	201
LINE	012	010.9	198
LTR	003	002.7	319
Simple_repeat	033	030.0	049
DNA	007	006.4	166
Low_complexity	040	036.4	047
snRNA	000	000.0	000
Other	000	000.0	000
tRNA	000	000.0	000
Satellite	000	000.0	000
rRNA	000	000.0	000
scRNA	000	000.0	000
Unknown	000	000.0	000
snpRNA	000	000.0	000
RNA	000	000.0	000
Total	110	100.0	

Control samples

- Well characterised case:
- e.g. Craig Venter (see: <http://ccr.coriell.org/sections/collections/HuRef/?SsId=78>). He's available for \$150 for 50ug to academic institutions.
- Settled for a CEPH trio that is included in the 1,000 genomes project

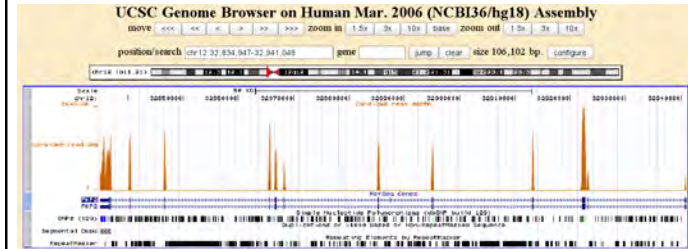


Exons analysed 1,325
 Total reads 19,875,187
 Reads not mapping to an exon 10,058823
 Cumulative exon length 446134 bases
 981 bases had zero coverage (as unique reads)
 1,874 had coverage <10

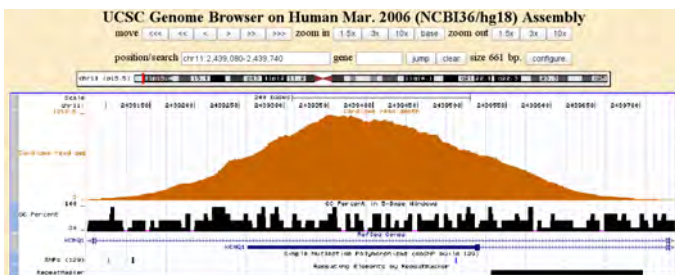
SNP Concordance

DBSNP	Target Array	Exome
identical	230	76
different	3	12
not found	0	89
not in CCDS		56
Total	233	233

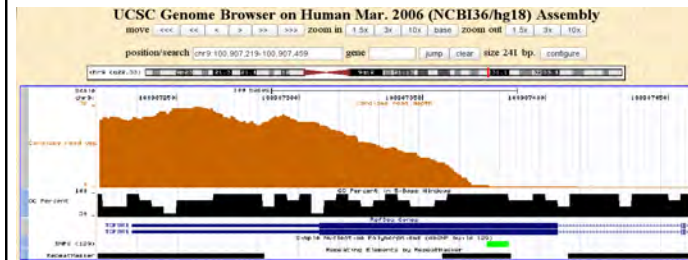
Coverage in a well-behaved region



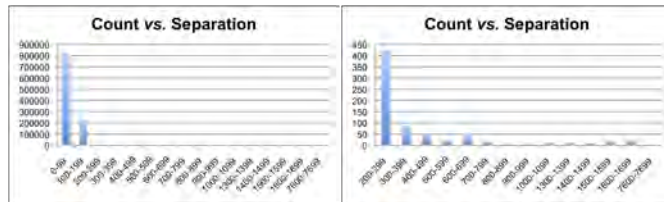
Good coverage extended by 150bp on either side



Inadequate coverage - why?

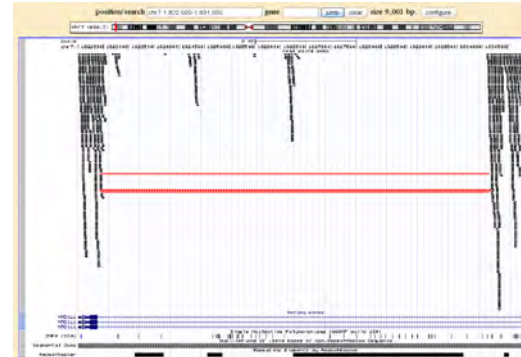


Genomic rearrangements by paired read analysis



Size of fragment sequenced was ~200b
 Separations > 200b may indicate large changes

A novel homozygous deletion?



Partitioning Summary

- A practical proposition
- Needs continuous review against whole genome sequencing (currently @ around £5,000)
- Some off target effects
- Some problems with Eland and multiple alignments (BWA is OK)
- Genome browsing and reporting still to be finalised

Copy number by sequence read depth

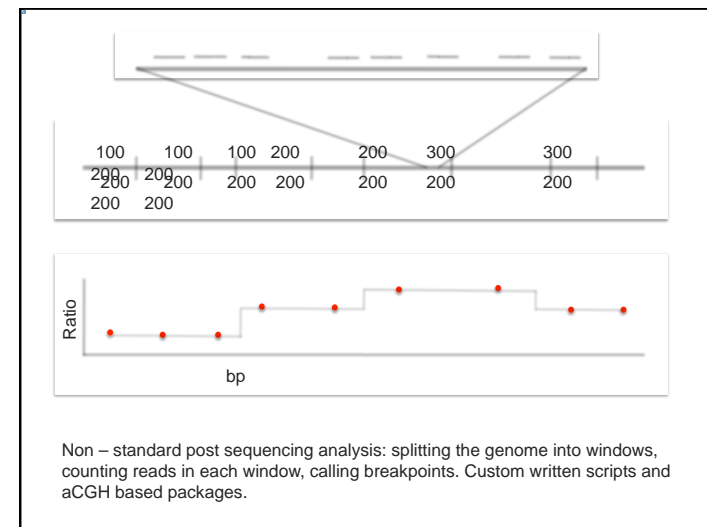
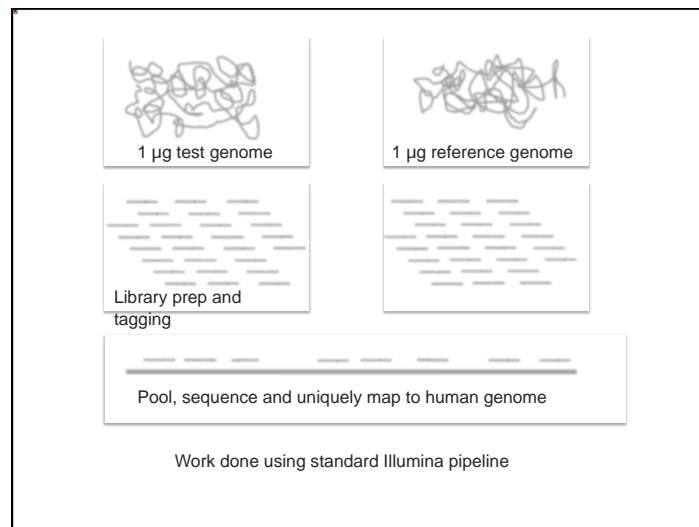
or – aCGH by next generation sequencing

Limitations of array CGH

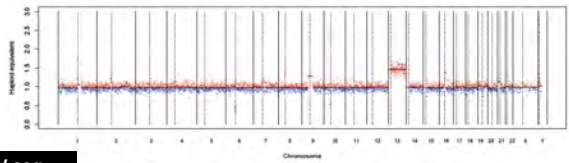
- Reliance on nucleic-acid hybridisation
 - DNA quality effects
 - Prior sequence knowledge required
 - “Analog” nature of fluorescence signal
- Reproducibility

CNV-Seq

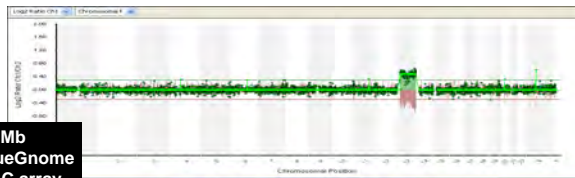
- Massively-parallel DNA sequencing using Illumina Genome II analyser
- ‘Tagging’ allows multiplex experiments per lane
- Standard and non-standard data analysis
- Sequence depth of coverage digitally quantified using custom-designed Python and R scripts



Aneuploidy detection

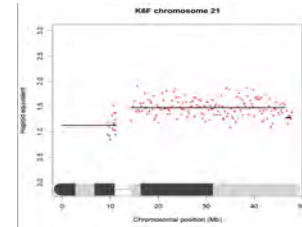


CNV-seq

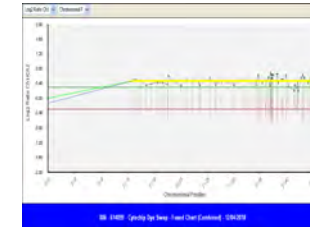


0.5Mb
BlueGenome
BAC array

Aneuploidy detection



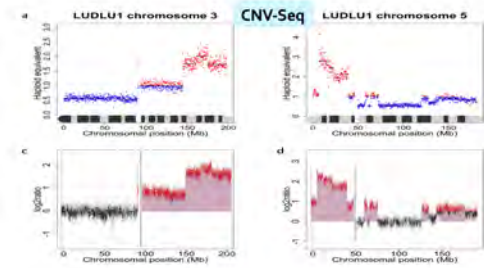
CNV-Seq



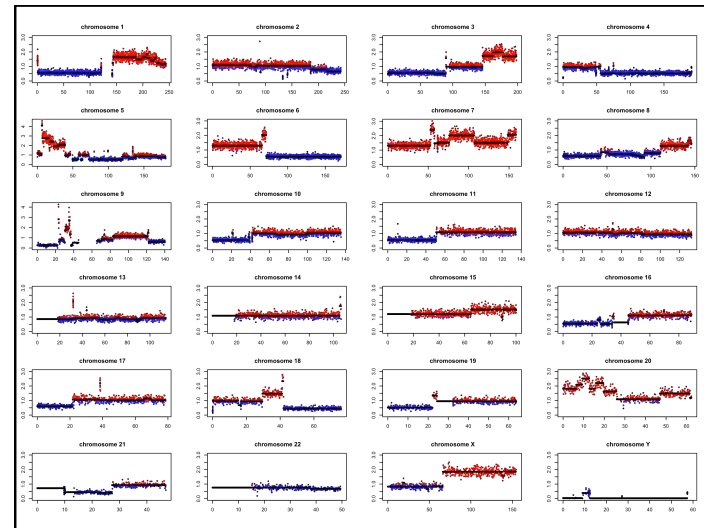
Array CGH

Chromosome 21 view showing trisomy

SCC Lung Tumour Cell lines



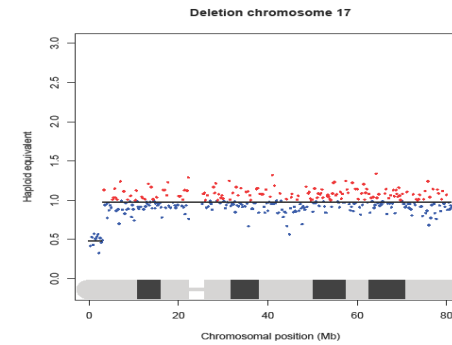
244K Agilent oligoarray



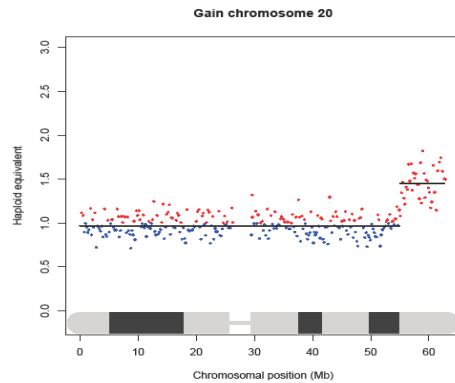
Prenatal CNV-Seq

- Bilateral ventriculomegaly on US at 21 weeks
- Amniocentesis
 - qfPCR normal
 - Long term culture normal karyotype
- CNV-Seq

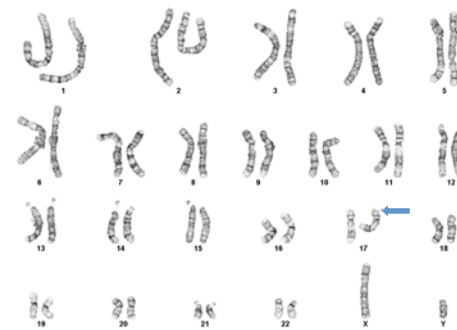
Chr 17: 1- 2 553 180 deletion (2.5Mb)

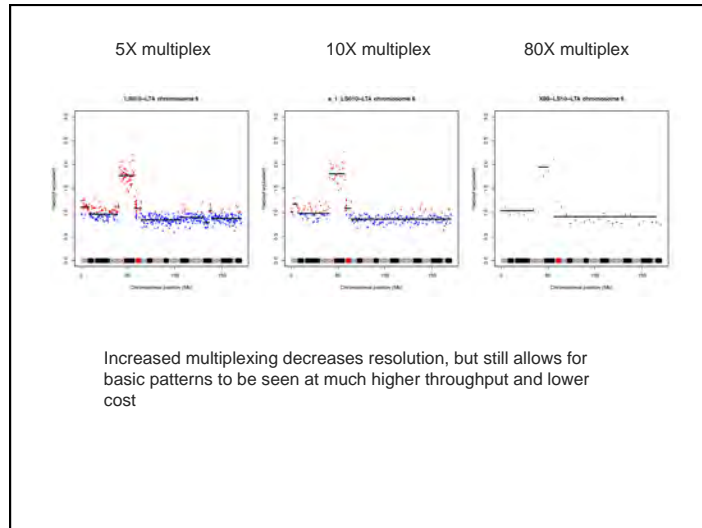


Chr 20: 54 963 376 – 62 332 309 (7.3 Mb gain)



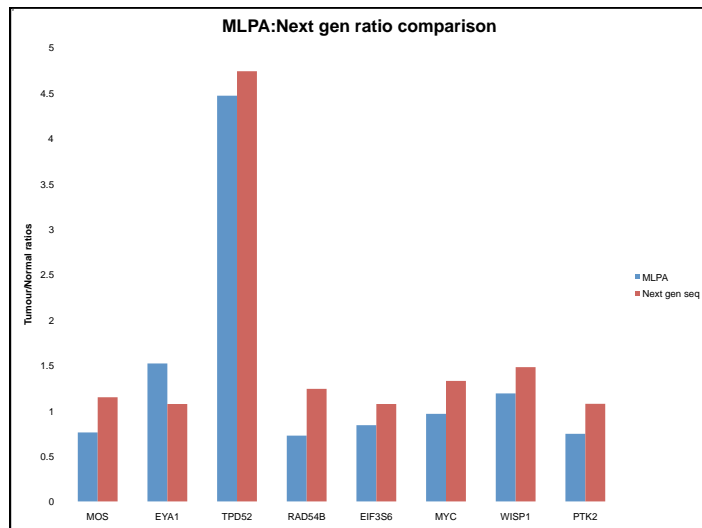
Fetal karyotype
46XY, der(17)t(17:20)(p13.3;q13.3) pat





Multiplexing and effect on resolution

Samples per lane (cost)	Number of reads	Mean distance between reads	Mean size of 200 read window
2 (£500)	5 million	620bp	123Kb
5 (£200)	2 million	1165bp	233Kb
10 (£100)	1 million	3062bp	612Kb
80 (£12)	125,000	27Kb	5.5Mb



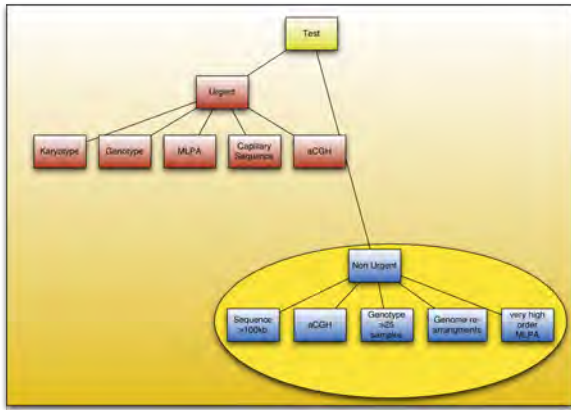
CNV-seq Summary

Next generation sequence platforms can produce copy number data from genomic DNA.

Data is comparable with aCGH data.

Tagging allows levels of multiplexing – hence throughput, cost and resolution can be tailored to the needs of a test.

Service design



Coming soon...

- Structural changes
- Circulating DNA
- Single molecule sequencing

Structural changes

Nature Genetics **40**, 722 - 729 (2008)
 Published online: 27 April 2008 | doi:10.1038/ng.128

Identification of somatically acquired rearrangements in cancer using genome-wide massively parallel paired-end sequencing

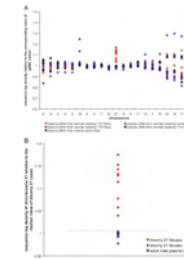
Peter J Campbell^{1,2,3}, Philip J Stephens^{1,2,3}, Erin D Pleasance¹, Sarah O'Meara¹, Heng Li¹, Thomas Santarius^{1,2,3}, Lucy A Stebbings¹, Catherine Leroy¹, Sarah Edkins¹, Claire Hardy¹, Jon W Teague¹, Andrew Menzies¹, Ian Goodhead¹, Daniel J Turner¹, Christopher M Clea¹, Michael A Quail¹, Antony Cox¹, Clive Brown¹, Richard Durbin¹, Matthew E Hurles¹, Paul A W Edwards², Graham R Bignell¹, Michael R Stratton¹ & P Andrew Futreal¹

Circulating DNA

Noninvasive diagnosis of fetal aneuploidy by shotgun sequencing DNA from maternal blood

H. Christina Fan¹, Yan J Blumenfeld², Usha Chikara¹, Lesanne Hudgins¹, and Stephen K Quake^{1*}

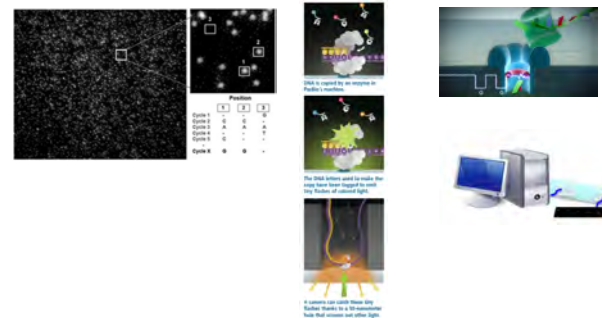
www.pnas.org/cgi/doi/10.1073/pnas.0808319105



Summary & Conclusions

- Clonal sequencing can be applied to diagnostic casework with high sensitivity and specificity
- The cost and efficiency benefits will depend on sequence quality, assay design and scaling
- The capacity is huge: 1 machine delivers over 20 Gigbases per week
- Clonal sequencing will contribute to improved genetic diagnostics

Single Molecule Sequencing



Acknowledgements

Funding

- Emmandjay Trust
- Cancer Research UK
- Dept Health NEAT
- LTH Challenge Fund

Science

- | | |
|------------------|---------------------|
| ◆ Jo Morgan | ◆ David Bonthron |
| ◆ Aengus Stewart | ◆ Colin Johnson |
| ◆ Heather Fraser | ◆ Ian Carr |
| ◆ Nick Camm | ◆ Constanze Bonifer |
| ◆ Helen Lindsay | ◆ Reuben Tooze |
| ◆ Chris Watson | ◆ Carol Chu |
| ◆ Kelly Cohen | ◆ Ruth Charlton |
| ◆ Bruce Hayward | ◆ Pamela Rabbitts |
| ◆ Henry Wood | ◆ Sir Alex Markham |
| ◆ Antigone Tzika | ◆ many others..... |