# Best-practice guidelines for reporting diagnostic NGS variants

Fran Smith

Clinical Scientist
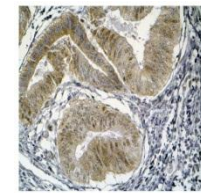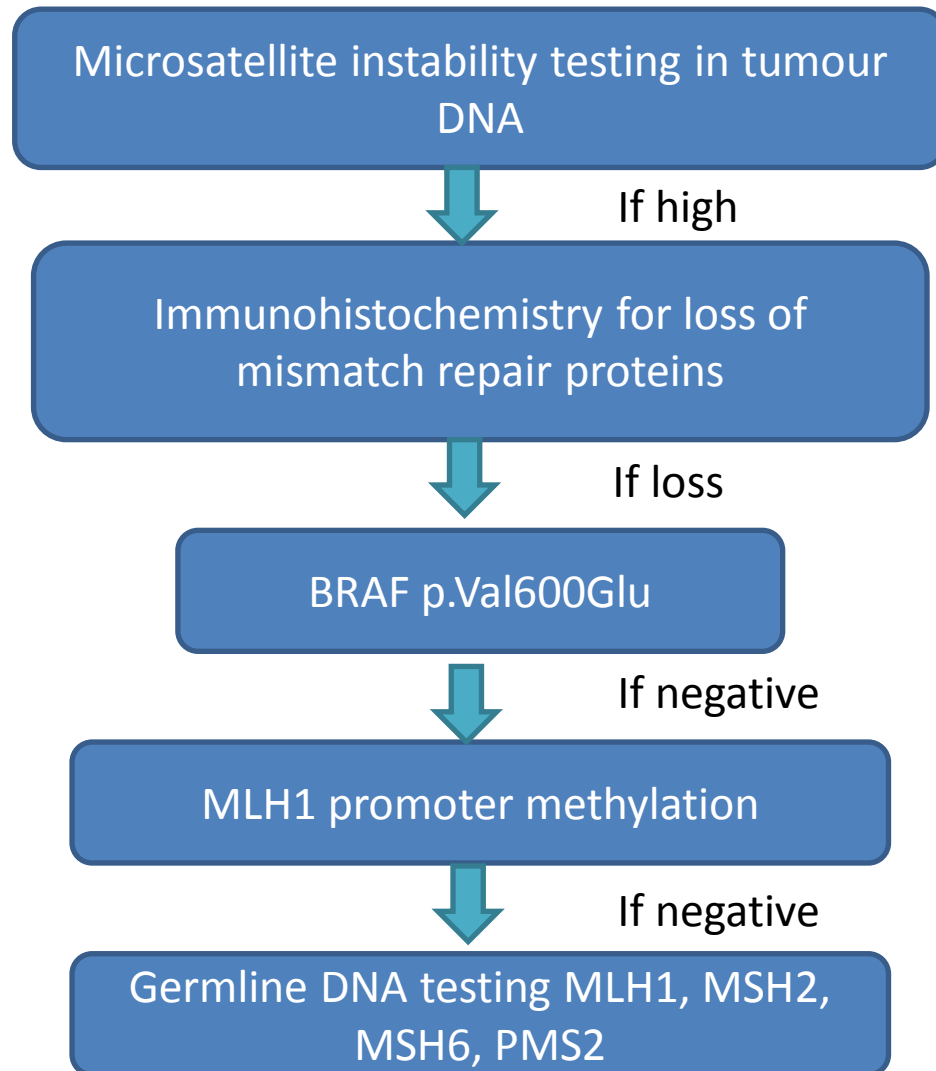
Molecular Pathology, King's College Hospital

# Overview:

- Background
  - Whole genome sequencing in context of clinical testing
- What we detect
- Variant classification systems and guidelines
  - ACGS
  - ACMG
- Terminology
- Functional effects
  - What significance do these variants have?
- Evidence we can use to interpret variants
  - Effect on protein
  - Splice site prediction tools
  - Missense variants
  - Databases
  - Frequency
  - Functional assays
  - Segregation
- Variant confirmation and validation
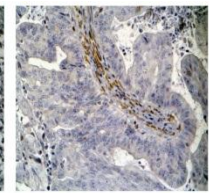  - What it means to be 'diagnostic'

# Context

- Genetics in diagnosis of disease has been very expensive
  - Pre 2012 £600 for one or two genes
  - Can now order a whole exome for about the same price
- Genetic testing was only undertaken when the clinician was fairly sure of the diagnosis
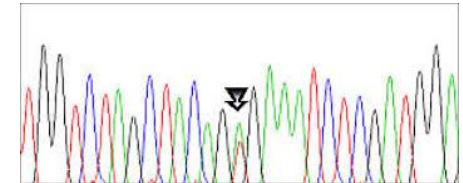- Now we can sequence all the genes and relate to phenotype after

# *Example: Lynch syndrome testing pathway*

Microsatellite instability testing in tumour DNA

↓ If high

Immunohistochemistry for loss of mismatch repair proteins

↓ If loss

BRAF p.Val600Glu

↓ If negative

MLH1 promoter methylation

↓ If negative

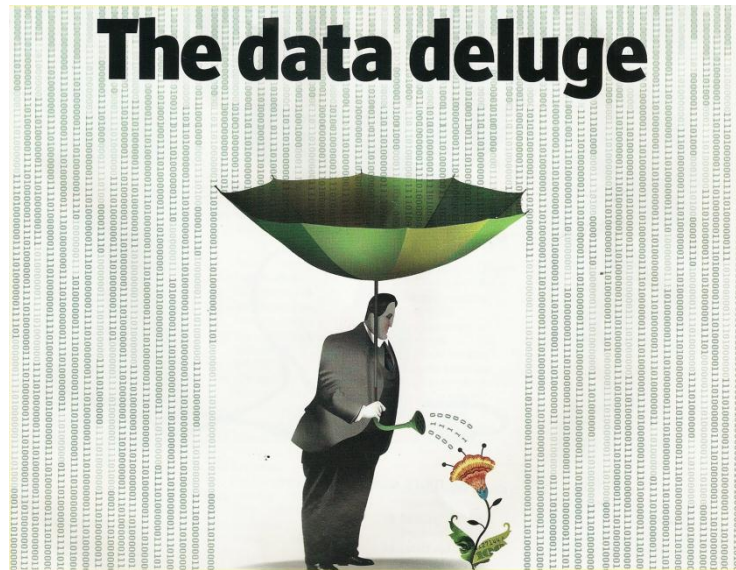Germline DNA testing MLH1, MSH2, MSH6, PMS2

# Whole genome sequencing (WGS)

- Only ~1 % of the genome codes for proteins
- Many other functional elements
- Numbers of variants expected from 1 genome:

| Type of variant | Average number in European | Average number in African |
|---|---|---|
| SNP | 3.53 million | 4.31 million |
| Indel | 546k | 625k |
| Large deletion | 939 | 1.1k |
| CNV | 157 | 170 |
| Non-synonymous | 10.2k | 12.2k |
| Synonymous | 11.2k | 13.8k |
| Intronic | 1.68 million | 2.06 million |

Data from 1000 genomes project

# Out of all these variants...which ones cause disease???



- We must name them systematically (HGVS)
- We must classify them using a standard evidence based approach

# HGVS nomenclature

- Must be called against a reference sequence
  - RefSeq
  - LRG – Locus Reference Genomic
  - Genome build
- g. is for genomic locations c. is for coding DNA sequence
- p. is for protein
- m. is for mitochondrial DNA
- Coding nomenclature is where c.1 is the first base of the translation initiation codon

# Transcripts

# Classification – UK & US system

| Class | Description | Report wording |
|-------|-------------|----------------|
| 1 | Clearly not pathogenic<br>Benign | Not reported |
| 2 | Unlikely to be pathogenic<br>Likely benign | Diagnosis not confirmed |
| 3 | Uncertain significance | Diagnosis not confirmed or excluded |
| 4 | Likely to be pathogenic | Consistent with diagnosis |
| 5 | Certainly pathogenic | Confirms diagnosis |

# Terms used to describe sequence variants

- **Pathogenic** - contributes mechanistically to disease, but is not necessarily fully penetrant

- **Damaging** - alters the normal levels or biochemical function of a gene or gene product – not necessarily causative

- **Deleterious** - reduces the reproductive fitness of carriers, targeted by purifying natural selection

- **Benign** – not harmful in effect

# Key considerations

- ## What is the expected inheritance pattern?
  - Autosomal dominant
  - Autosomal recessive
  - X-linked
  - De novo

- ## What is the mechanism of pathogenicity?
  - Loss of function
  - Gain of function
  - Haploinsufficiency

- ## What is the frequency of the disease?
  - Can you rule out common variants (>1%)?

# Evidence base for classification

- ACMG (American College of Medical Genetics) have very useful descriptive classification system:

| Pathogenic | | Example |
|---|---|---|
| PVS1 | Very Strong | Nonsense variant |
| PS1-6 | Strong | Well established functional assay |
| PM1-6 | Moderate | In mutational hotspot |
| PP1-5 | Supporting | Cosegregates with disease |

| Benign | | Example |
|---|---|---|
| BA1 | Stand-alone | Allele frequency >5% |
| BS1-4 | Strong | Lack of effect on functional assay |
| BP1-7 | Supporting | Observed in *trans* with pathogenic variant |

**Figure 1 Evidence framework.**

| | Benign | | Pathogenic | | | |
| Evidence type | Strong | Supporting | Supporting | Moderate | Strong | Very strong |
|---|---|---|---|---|---|---|
| **Population data** | MAF is too high for disorder BA1/BS1 **OR** observation in controls inconsistent with disease penetrance BS2 | | | Absent in population databases PM2 | Prevalence in affecteds statistically increased over controls PS4 | |
| **Computational and predictive data** | | Multiple lines of computational evidence suggest no impact on gene /gene product BP4<br><br>Missense in gene where only truncating cause disease BP1<br><br>Silent variant with non predicted splice impact BP7<br><br>In-frame indels in repeat w/out known function BP3 | Multiple lines of computational evidence support a deleterious effect on the gene /gene product PP3 | Novel missense change at an amino acid residue where a different pathogenic missense change has been seen before PM5<br><br>Protein length changing variant PM4 | Same amino acid change as an established pathogenic variant PS1 | Predicted null variant in a gene where LOF is a known mechanism of disease PVS1 |
| **Functional data** | Well-established functional studies show no deleterious effect BS3 | | Missense in gene with low rate of benign missense variants and path. missenses common PP2 | Mutational hot spot or well-studied functional domain without benign variation PM1 | Well-established functional studies show a deleterious effect PS3 | |
| **Segregation data** | Nonsegregation with disease BS4 | | Cosegregation with disease in multiple affected family members PP1 | → Increased segregation data → | | |
| **De novo data** | | | | De novo (without paternity & maternity confirmed) PM6 | De novo (paternity and maternity confirmed) PS2 | |
| **Allelic data** | | Observed in *trans* with a dominant variant BP2<br><br>Observed in *cis* with a pathogenic variant BP2 | | For recessive disorders, detected in trans with a pathogenic variant PM3 | | |
| **Other database** | | Reputable source w/out shared data = benign BP6 | Reputable source = pathogenic PP5 | | | |
| **Other data** | | Found in case with an alternate cause BP5 | Patient's phenotype or FH highly specific for gene PP4 | | | |

This chart organizes each of the criteria by the type of evidence as well as the strength of the criteria for a benign (left side) or pathogenic (right side) assertion. Evidence code descriptions can be found in **Tables 3** and **4**. BS, benign strong; BP, benign supporting; FH, family history; LOF, loss of function; MAF, minor allele frequency; path., pathogenic; PM, pathogenic moderate; PP, pathogenic supporting; PS, pathogenic strong; PVS, pathogenic very strong.

# What effect does the variant have?

- Is it in a region known to be functional?
  - Protein coding regions
  - Flanking splice site
  - Promoter
- However we know from ENCODE that there are many more regions of the genome that are functional

In a known functional region = more likely to be pathogenic

# Types of variant that will severely affect the protein ...*there are always exceptions*

- Nonsense
  - Premature stop codons – nonsense mediated decay (NMD)
- Frameshift
  - Translated protein scrambled – often results in NMD
- Canonical splice site
  - Causes exon skipping and loss of functional domains or frameshift

# Nonsense

Normal sequence

DNA sequence

| ATG | GGA | AGA | CCG | TCC | TGA |
|-----|-----|-----|-----|-----|-----|
| Met | Gly | Arg | Pro | Ser | * |

Protein sequence

Mutated sequence          c.7A>T

DNA sequence

| ATG | GGA | TGA |
|-----|-----|-----|
| Met | Gly | * |

Protein sequence

Loss of part of functional part of protein = likely to be pathogenic

# Exception to the rule…

- A truncating variant in the last exon
- BRCA2: c.9976A > T (p.Lys3326*)
- Protein is still functional



- Common in European population
- Was mistakenly assigned as a pathogenic variant
- Always be aware of the context of your mutation within the gene

# Frameshift

Scrambled reading frame = likely to be pathogenic

Normal sequence

DNA sequence

Protein sequence

| ATG | GGA | AGA | CCG | TCC | TGA |
|-----|-----|-----|-----|-----|-----|
| Met | Gly | Arg | Pro | Ser | *   |

Mutated sequence          c.7delT

DNA sequence

Protein sequence

| ATG | GGA | GAC | CGT | CCT | GA... |
|-----|-----|-----|-----|-----|-------|
| Met | Gly | Asp | Arg | Pro | ..... |

# Splice site

- Splicing is all about recognition of sequence motifs by the spliceosome:



- Variants in the consensus splice site can lead to exon skipping or incorporation of intronic sequence into the mRNA transcript
- Cryptic splice sites can also be activated

# Splice site prediction software

# Subtle effects on protein - missense

- Substituting one amino acid for another
- Can affect functional part of protein
- Adding or removing cysteine alters potential for forming disulphide bridges
- Grantham distance: a measure of physiochemical difference

Large physiochemical difference = more likely to be pathogenic

# Amino acid properties

# Grantham distance

# Multiple sequence alignments



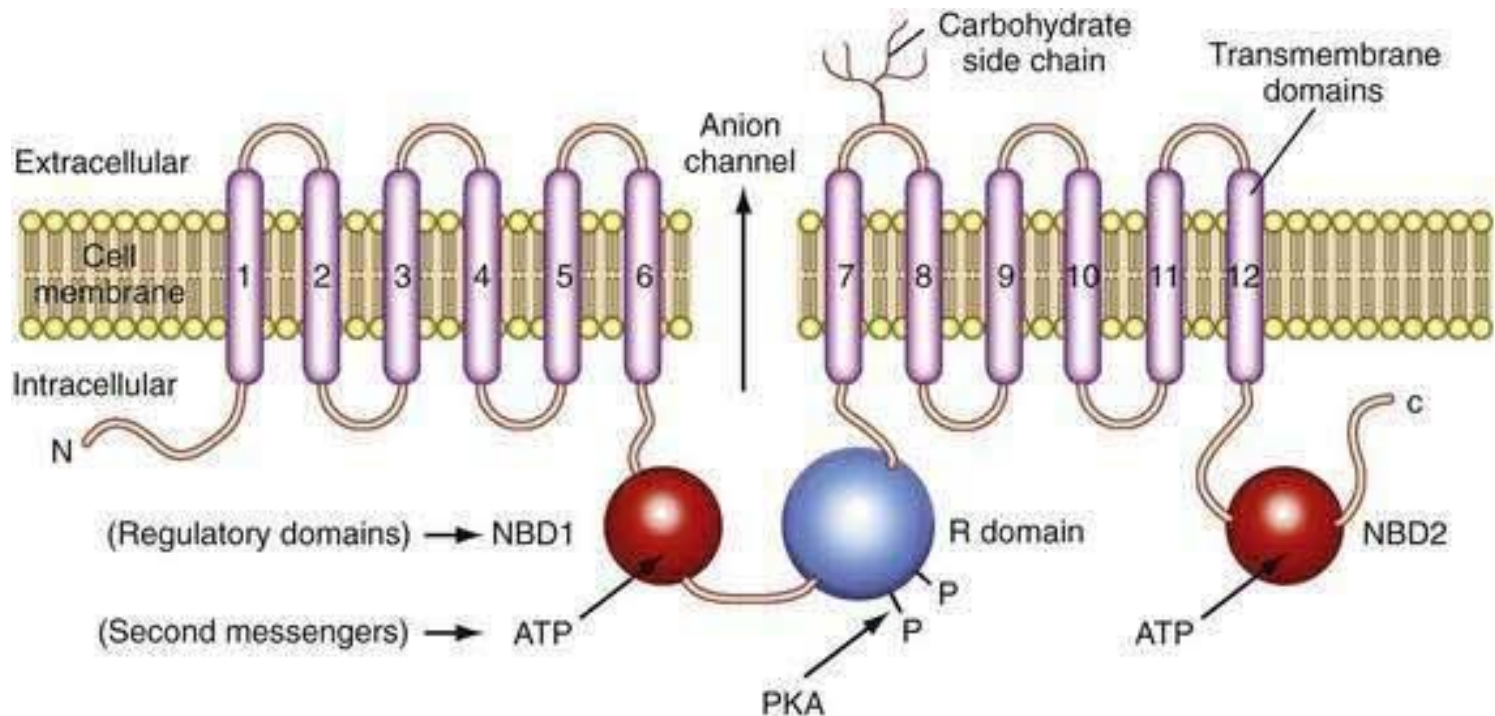- Assessment of evolutionary conservation

More conserved = more important for function

Not conserved = less likely to be pathogenic

# Protein domains/mutational hotspots

CFTR protein structure



In important protein domain = more likely to be pathogenic

# Missense prediction tools

- In silico prediction
- SIFT – Sorting intolerant from tolerant
- PolyPhen
- AlignGVGD
- MutationTaster
- All are based upon physiochemical differences, position in the protein and multiple sequence alignments or a combination
- Use with CAUTION
- Never use in isolation

Very difficult to tell anything from missense prediction tools!

# HBB c.20A>T; p.Glu7Val
# HbS Sickle

- One example exemplifies why you must consider ALL THE DATA
- Missense
- 4.85% in African population (ExAC)
- Weakly conserved nucleotide (phyloP: 0.04 [-14.1;6.4])
- Moderately conserved amino acid (considering 20 species)
- Moderate physicochemical difference between Glu and Val (Grantham dist.: 121 [0-215])
- This variant is in protein domain: Globin, subset
- Align GVGD: C0 (GV: 164.97 - GD: 0.00)
- SIFT: Tolerated (score: 0.1, median: 2.21)
- MutationTaster: polymorphism (p-value: 1)

# Databases

- Leiden Open Variation database LOVD
- Human Gene Mutation Database HGMD
- ClinVar
- Online Mendelian Inheritance in Man OMIM
- DECIPHER
- Many smaller disease specific databases
- Considerations
  - Is it curated/updated?
  - HGVS nomenclature and transcript
  - Validated data
  - Source and independence of the observations listed
  - Always reassess their data

# Variant frequency

Higher frequency = less likely to be pathogenic

- Probably one of the most useful ways of assigning classes 1 and 2 (not pathogenic)
- Consider the incidence of the disease vs variant frequency and inheritance
- Several large databases that hold variant frequency data
  - ExAC – Exome Aggregation consortium
  - dbSNP
  - EVS – Exome variant server
- According to ACMG: >5% freq is stand-alone support for classification as benign

# Example:

- You find a variant at 36% in the Latino population in ExAC

**Population Frequencies**

| Population | Allele Count | Allele Number | Number of Homozygotes | Allele Frequency |
|---|---|---|---|---|
| East Asian | 3594 | 8570 | 748 | 0.4194 |
| Latino | 4183 | 11494 | 794 | 0.3639 |
| European (Finnish) | 2285 | 6604 | 391 | 0.346 |
| Other | 275 | 900 | 46 | 0.3056 |
| European (Non-Finnish) | 18446 | 66578 | 2563 | 0.2771 |
| South Asian | 2936 | 16506 | 305 | 0.1779 |
| African | 856 | 9792 | 32 | 0.08742 |
| Total | 32575 | 120444 | 4879 | 0.2705 |

- What are your considerations?
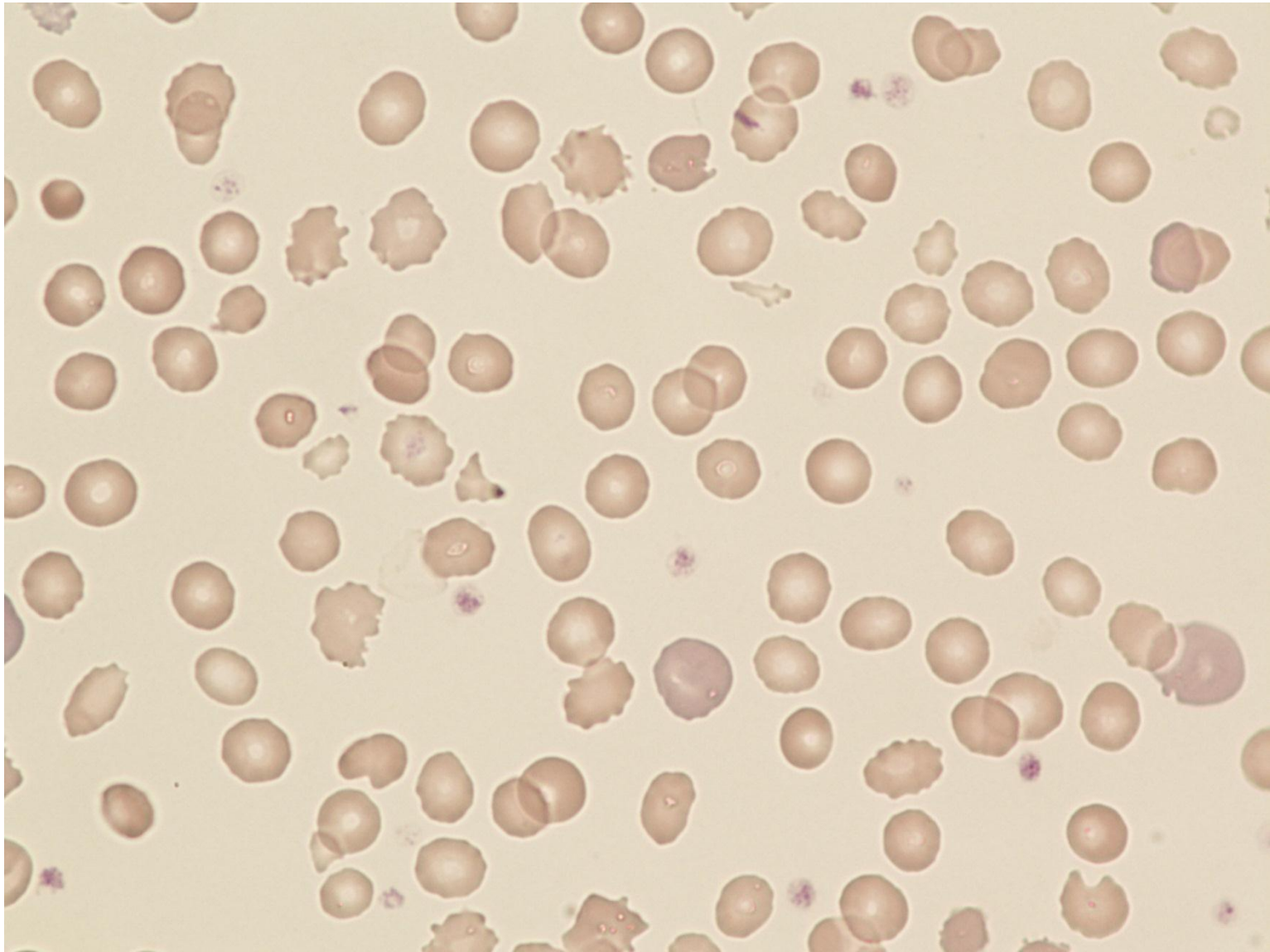
# Functional assays

- Can be very useful in looking at the *in vivo* effect of a variant

- Many metabolic enzymes have robust functional assays
  - Pyruvate kinase

- mRNA sequencing very good at predicting effects of splicing variants

# Example:

- You find a PKLR class 3 missense in *trans* with a class 5 in a patient with haemolytic anaemia

- PKLR activity can be assayed on a fresh EDTA sample

- PKLR assay shows the patient to be deficient

- This data can help to re-classify the class 3 variant as pathogenic

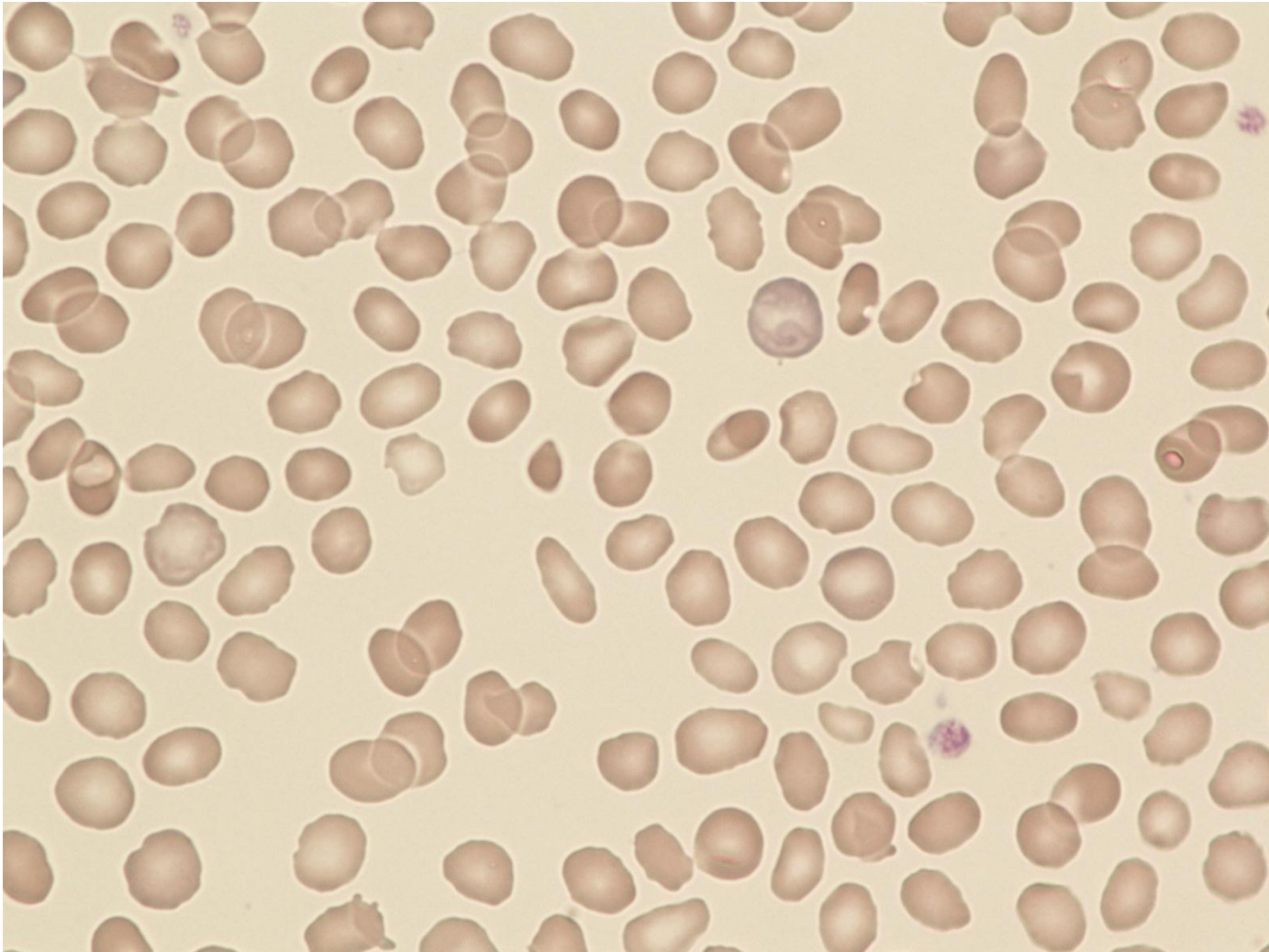- Blood film: You are looking at the functional effect

- Mother:
  - Homozygous *SPTA1:* c.[5572C>G; 6531-12C>T]; p.[Leu1858Val;?] Low expression allele

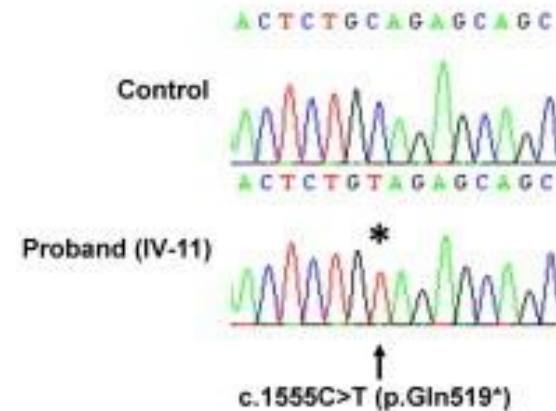- Father:
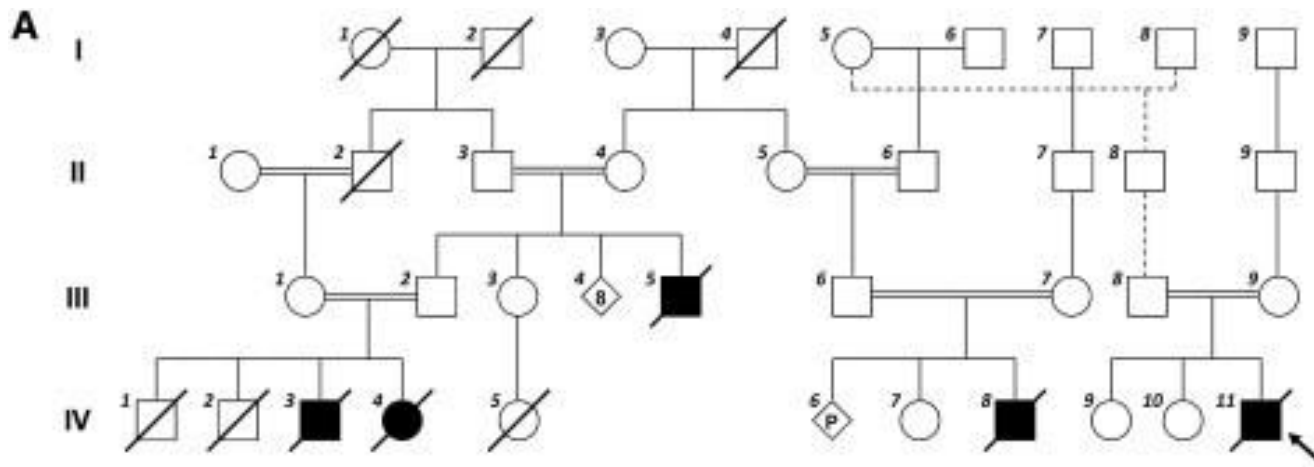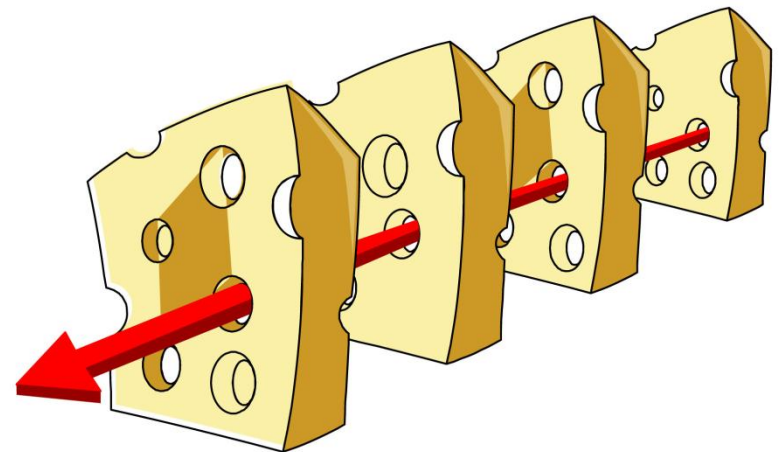  - Heterozygous *SPTA1:* c.83G>A; p.Arg28His

# Segregation

- Is it in *cis/trans* with other variants in the patient?
- Testing of other affected or unaffected family members

# Diagnostic validation of results

- What does it mean to be 'diagnostic'?
  - ISO15189 accreditation for quality standards
  - The cheese model
- We confirm a variant by Sanger sequencing
- Reassess pathogenicity in relation to phenotype
- Report

# What should be included on the clinical report?

- Results
  - HGVS nomenclature
  - Gene name, cDNA, protein
  - Disease
  - Inheritance
- Interpretation
  - Evidence supporting variant classification
  - Does it explain the patient's phenotype?
  - Supplementary testing

# Report…

- Methodology
  - Laboratory and analysis tools used
  - Limitations
- Risk to offspring
- Referral for genetic counselling
- Incidental findings??

# Incidental findings

- Much debate on if we should report pathogenic variants that are not associated with the condition that the patient has been referred for
  - Eg Cystic Fibrosis carrier status
  - Eg Late onset dominant disorders such as BRCA
- For 100k the patient will decide at consent

# Examples

- Hereditary spastic paraplegia referral
- Heterozygous *REEP1* variant c.471del p.(Thr158fs)
- What questions will you ask?

- Is this gene known to be associated with the condition?
  - Yes
- Mode of inheritance?
  - Autosomal dominant
- How big is the gene?
  - 7 exons
- Where does the variant lie in the gene
  - Penultimate exon