



LEIDEN UNIVERSITY MEDICAL CENTER

# **Diagnostic Variant Database**

## **an overview**

**Jeroen F. J. Laros**

**Leiden Genome Technology Center**

**Department of Human Genetics**

**Center for Human and Clinical Genetics**



*Share your results*

*Share your results*

There are a number of ways of storing and sharing your variants:

- LOVD.
- dbSNP.
- ...


## *Share your results*

There are a number of ways of storing and sharing your variants:


- LOVD.
- dbSNP.
- ...

Sharing is important:

- If no-one shared, you can not filter.
- Can lead to co-authorship when other people use your data.



## Leiden Muscular Dystrophy pages

Duchenne Muscular Dystrophy (DMD) 

Curator: [Johan den Dunnen](#)

LOVD v.2.0 Build 34 [ [Current LOVD status](#) ]  
[Register as submitter](#) | [Log out](#)

Home
Variants
Submitters
Submit
Documentation

DMD homepage [Switch gene](#)

When referring to this database please cite [Aartsma-Rus et al. \(2006\). Entries in the Leiden Duchenne muscular dystrophy mutation database: an overview of mutation types and paradoxical cases that confirm the reading-frame rule. Muscle Nerve. 34:135-144.](#)

**NOTE: for deletions / duplications go to the [DMD database for whole exon changes](#).**

### LOVD Gene homepage

General information	
Gene name	Duchenne Muscular Dystrophy
Gene symbol	<b>DMD</b>
Chromosome Location	Xp21.2
Database location	<a href="#">the Leiden Muscular Dystrophy pages</a>
Curator	<a href="#">Johan den Dunnen</a>
PubMed references	View all (unique) <a href="#">PubMed references</a> in the DMD database
Date of creation	July 29, 1997
Last update	September 21, 2012
Version	<b>DMD120921</b>
Add sequence variant	<a href="#">Submit a sequence variant</a>
First time submitters	<a href="#">Register here</a>
Reference sequence file	<a href="#">coding DNA reference sequence</a> for describing sequence variants
Genomic refseq ID	<a href="#">NG_012232.1</a>
Transcript refseq ID	<a href="#">NM_004006.2</a>
Exon/intron information	<a href="#">Exon/intron information table</a>

Figure 1: LOVD welcome screen.

## Variants in LOVD

Exon	DNA change	Var_pub_as	RNA change	Protein change
6	c.368dupC	-	r.(?)	p.(Met124Asnfs*14)
6	c.372delG	-	r.372del	p.Met124Ilefs*18
6	c.377delA	-	r.(?)	p.(Asn126Ilefs*16)
6	c.379_380del (Reported 4 times)	-	r.379_380del	p.Ile127Hisfs*10
6	c.386delC	-	r.(?)	p.(Ala129Valfs*13)
6	c.392T>A	-	r.(?)	p.(0)
6	c.397C>T	-	r.(?)	p.(Gln133*)
6	c.402_405del (Reported 2 times)	401_404delCCAA	r.(?)	p.(Asn135Valfs*6)
6	c.409G>T	-	r.(?)	p.(Glu137*)
6	c.419T>A (Reported 2 times)	-	r.(?)	p.(Leu140His)
6	c.429G>A	-	r.(?)	p.(Trp143*)
6	c.431T>A (Reported 2 times)	-	r.(?)	p.(Val144Asp)
6	c.433C>T (Reported 19 times)	-	r.(?)	p.(Arg145*)
6	c.433delC	-	r.(?)	p.(Arg145Aspfs*12)
6	c.440C>A	-	r.440c>a	p.Ser147*
6	c.440C>G (Reported 2 times)	-	r.(?)	p.(Ser147*)

Figure 2: Selection of variants.

*DVD*

The *Diagnostic Variant Database*.

- Share variants found in exome sequencing experiments.
- Find functionally relevant variants.

## *DVD*

The *Diagnostic Variant Database*.

- Share variants found in exome sequencing experiments.
- Find functionally relevant variants.

Technical details:

- Annotation by sharing.
- Store coverage information to determine reference calls.
- Disambiguation of variant descriptions.
- Pooling without loss of information.
- Duplicate sample detection.
  - Allows for re-annotation without polluting the database.
- Encrypted connection with authentication.



## *Data generation*



Figure 3: SOLiD 5500XL.



Figure 4: HiSeq 2000.

## Data analysis

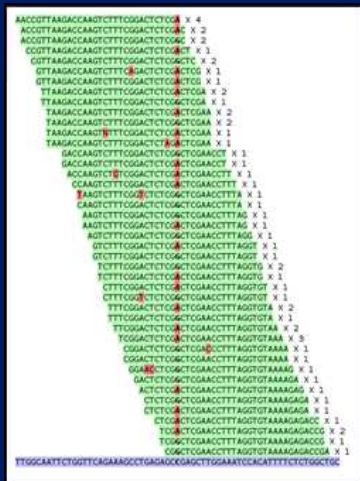


Figure 5: Variant calling.

### *Required fields in the VCF format*

Field	Explanation.
CHROM	Name of the chromosome.
POS	Position on the chromosome.
ID	List of unique identifiers.
REF	Reference base(s).
ALT	List of alternate non-reference alleles.
QUAL	Phred-scaled quality score for the assertion made in ALT.
FILTER	PASS if this position has passed all filters.
INFO	Additional information.

Table 1: Required fields.

## *Disambiguation*

One problem with the VCF format is that there are multiple ways of describing the *same* variant.

Original	Stored
<b>AT</b> ⇒ <b>A</b>	<b>AT</b> ⇒ <b>A</b>
<b>ATT</b> ⇒ <b>AT</b>	<b>AT</b> ⇒ <b>A</b>
<b>ATTC</b> ⇒ <b>ATC</b>	<b>AT</b> ⇒ <b>A</b>
<b>G</b> ⇒ <b>AG</b>	<b>G</b> ⇒ <b>AG</b>
<b>GG</b> ⇒ <b>AGG</b>	<b>G</b> ⇒ <b>AG</b>
<b>GGA</b> ⇒ <b>AGGA</b>	<b>G</b> ⇒ <b>AG</b>

Table 2: Resolving ambiguous variant descriptions.

## *Disambiguation*

One problem with the VCF format is that there are multiple ways of describing the *same* variant.

Original	Stored
<b>AT</b> ⇒ <b>A</b>	<b>AT</b> ⇒ <b>A</b>
<b>ATT</b> ⇒ <b>AT</b>	<b>AT</b> ⇒ <b>A</b>
<b>ATTC</b> ⇒ <b>ATC</b>	<b>AT</b> ⇒ <b>A</b>
<b>G</b> ⇒ <b>AG</b>	<b>G</b> ⇒ <b>AG</b>
<b>GG</b> ⇒ <b>AGG</b>	<b>G</b> ⇒ <b>AG</b>
<b>GGA</b> ⇒ <b>AGGA</b>	<b>G</b> ⇒ <b>AG</b>

Table 2: Resolving ambiguous variant descriptions.

This is done automatically on the server.

## *Coverage*

We need to store coverage information.

- Which variants *could* we have seen?

## *Coverage*

We need to store coverage information.

- Which variants *could* we have seen?

Frequency calculation:

- A variant can be present only twice in the database.
- If it is only covered in two samples, the frequency is one.

## *Coverage*

We need to store coverage information.

- Which variants *could* we have seen?

Frequency calculation:

- A variant can be present only twice in the database.
- If it is only covered in two samples, the frequency is one.

This enables us to store multiple experiment types.

- Different capture kits.
- Full genome sequencing.
- Amplicon resequencing.



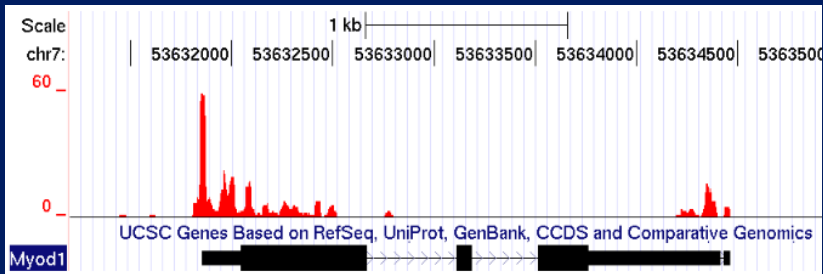
*Coverage*

Figure 6: Coverage / depth histogram.

## Coverage

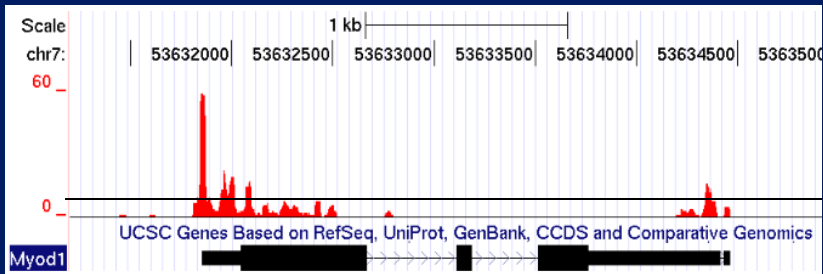


Figure 6: Coverage / depth histogram.

## Coverage

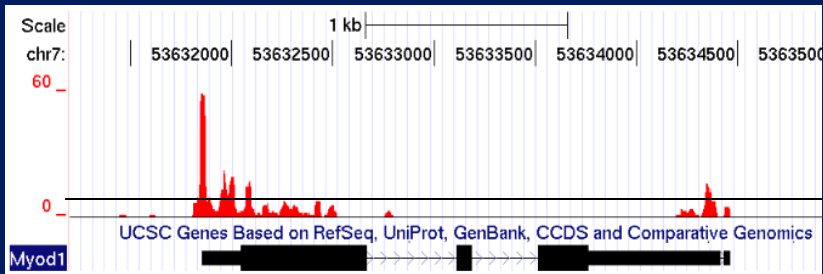


Figure 6: Coverage / depth histogram.



Figure 7: Coverage summary.

## *Minimal BED format*

We only store minimal information about coverage.

Field	Explanation.
CHROM	Name of the chromosome.
BEGIN	First position of the region.
END	Last position of the region.

Table 3: Three column BED format.

## *Pooling*

We provide a script to merge multiple VCF files into one.

- VCF files are merged and sorted on the *client* side.
  - The server has no knowledge of the individual variant calls.
- This makes identification of a person hard if the pool size is large enough.

## *Pooling*

We provide a script to merge multiple VCF files into one.

- VCF files are merged and sorted on the *client* side.
  - The server has no knowledge of the individual variant calls.
- This makes identification of a person hard if the pool size is large enough.

As long as we calculate the *coverage* files separately, we can still calculate the frequencies correctly.

- One variant file and multiple coverage files are sent to the server.

## *Pooling*

We provide a script to merge multiple VCF files into one.

- VCF files are merged and sorted on the *client* side.
  - The server has no knowledge of the individual variant calls.
- This makes identification of a person hard if the pool size is large enough.

As long as we calculate the *coverage* files separately, we can still calculate the frequencies correctly.

- One variant file and multiple coverage files are sent to the server.

Anonymity with no loss of functionality.

## Database layout

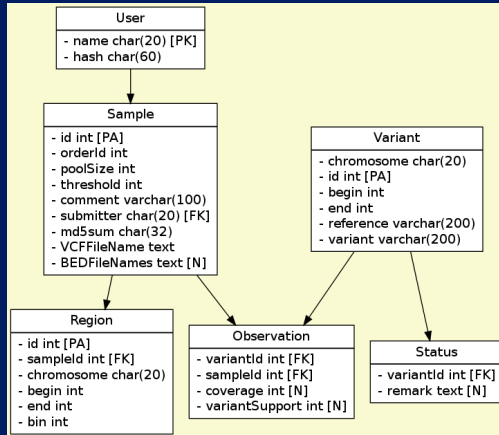


Figure 8: DVD SQL schema.



## *Set up*

For privacy reasons:

- No browsing.
- Pre-defined interface.
- Encrypted connection.
- Authentication.

## *Set up*

For privacy reasons:

- No browsing.
- Pre-defined interface.
- Encrypted connection.
- Authentication.

Protocols:

- A SOAP interface for interaction with pipelines.
- SOAP over HTTPS (SSL) for encryption.
- Authentication.
  - Send username / password in each call.
  - Only the password hash is stored in the database.

## *Client perspective*

In principle, one program that handles all:

- Upload a VCF and one or more BED tracks.
- Server responds with an ID.
- Use the ID to poll for the results.

## *Client perspective*

In principle, one program that handles all:

- Upload a VCF and one or more BED tracks.
- Server responds with an ID.
- Use the ID to poll for the results.

Client characteristics:

- Around 30 lines of code.
- Programming language independent.
- Operating system independent.

## *Server perspective*

### Uploading:

- Receive a VCF and one or more BED tracks.
- Respond with an ID.
- A process is forked that starts annotation and importing.

*Server perspective*

## Uploading:

- Receive a VCF and one or more BED tracks.
- Respond with an ID.
- A process is forked that starts annotation and importing.

## Annotation:

- Receive an ID.
- If the annotation for this ID is finished, return the result.

## *Server perspective*

### Uploading:

- Receive a VCF and one or more BED tracks.
- Respond with an ID.
- A process is forked that starts annotation and importing.

### Annotation:

- Receive an ID.
- If the annotation for this ID is finished, return the result.

Note that once the annotation part is finished, the result will be returned (before the import needs to be finished).

## *Miscellaneous*

The uploaded files are cached and *md5sums* are stored in the database.

- Prevents re-uploading the same data.
- Allows for re-annotation without polluting the database.
- Can be used to re-populate the database if needed.



## *Miscellaneous*

The uploaded files are cached and *md5sums* are stored in the database.

- Prevents re-uploading the same data.
- Allows for re-annotation without polluting the database.
- Can be used to re-populate the database if needed.

Highly optimised for genomic data.

- Binning to optimise region (spacial) data.
- Thoroughly tested and benchmarked.

*Output*

Column	Name	Description
1	CHROM	Name of the chromosome.
2	POS	Position on the chromosome.
3	REF	Reference allele.
4	ALT	Genotype.
5	Frequency	Frequency of this genotype.
6	Occurrences	Number of hits.
7	IDs (> 100)	If the frequency is low: a list of sample IDs.
8	IDs (< 100)	A list with hits in reserved samples; 1000 genomes, GoNL, HGMD, etc.

Table 4: The output format.

## *Production*

Currently setting up a server at the UMCG.

- Nijmegen.
- Utrecht.
- Leiden.
- Amsterdam.
- Groningen.

<https://trac.nbic.nl/dvd/>

## *Production*

Currently setting up a server at the UMCG.

- Nijmegen.
- Utrecht.
- Leiden.
- Amsterdam.
- Groningen.

Open source.

The core can be modified easily.

- Drop restrictions.
- Interface for browsing.
- ...

<https://trac.nbic.nl/dvd/>

## Acknowledgements:

Martijn Vermaat  
Christian Gilissen  
Victor Guryev  
Frans-Paul Ruzius  
Nienke Wieskamp  
Wim Spee  
Morris Swertz  
Pieter Neerincx  
Rob Hooft  
David van Enckevort  
Hailiang Mei  
Johan den Dunnen

<https://trac.nbic.nl/dvd/>