

# **Using Machine Learning and Expert Human Guidance to Automate Clinical Data Integration to CDISC Standards**

Timothy Danford, Tamr Inc.

# **Files, Names, & Standards**

Please Don't Make a  
Programmer Cry

# File Formats are Complicated

[illegible]

# Format Specification is Hard

```
/*-----*/
/* Name: ieee2xpt                                     */
/* Purpose: converts IEEE to transport                */
/* Usage: rc = ieee2xpt(to_ieee,p_data);              */
/* Notes: this routine is an adaptation of the wzctdbl routine */
/* from the Apollo.                                   */
/*-----*/
```

```
void ieee2xpt(ieee,xport)
```

```
unsigned char *ieee; /* ptr to IEEE field (2-8 bytes) */
unsigned char *xport; /* ptr to xport format (8 bytes) */
{
```

```
register int shift;
unsigned char misschar;
int ieee_exp;
unsigned long xport1,xport2;
unsigned long ieee1 = 0;
unsigned long ieee2 = 0;
```

# Are Your Formats Self-Describing?

- Are your formats even defined at all?
  - Programmers have a lot of tools to describe formats, they're called **grammars & semantics**
- Where is the format "flexible?"
  - Your data is not a **special snowflake**.
- Can the format describe itself?
  - How much communication is **out-of-band** (i.e. outside the format itself)?

# Naming Is Complicated

Reference and alternative alleles of a multi nucleotide polymorphism (MNP)

REF  
ALT

GGGCATGGG  
GGGTGCGGG

Genome Reference		Variant Call Format			
GGGGCATGGGG		POS	REF	ALT	
REF	GCAT	4	GCAT	GTGC	Not left trimmed
ALT	GTGC				
REF	CATG	5	CATG	TGCG	Not right trimmed
ALT	TGCG				
REF	GCATG	4	GCATG	GTGCG	Not left and right trimmed
ALT	GTGCG				
REF	CAT	5	CAT	TGC	Normalized
ALT	TGC				
Alleles represented against the human genome reference. Allele pairs are colored the same, all are representations of the same variant.		Alleles represented in Variant Call Format, all are representations of the same variant.			



# **Study Data Tabulation Model Implementation Guide: Pharmacogenomics/Genetics Version 1.0 (Provisional)**

**Prepared by the  
CDISC PGx Team**

## **Notes to Readers**

- This is the provisional version 1.0 of the implementation guide for Pharmacogenomics and Pharmacogenetics. It is intended to correspond to version 1.5 of the CDISC Study Data Tabulation Model.
- Because SDTM v1.5 has not yet been published, this document remains a provisional release only.

# “Multiple Names in Your Standard” Means You Have Multiple Standards

## 6. Genetic Variation:

- a. PFTESTCD and PFTEST generally identify the type of test performed by specifying the type of material assessed, such as nucleotides or amino acid, or the level of granularity, such as codon or allele. See the table below for suggested PFTESTCD values:

PFTESTCD	PFTEST	Notes
NUC	Nucleotide	Observes nucleotide sequences or values. Performed on DNA or RNA.
CDN	Codon	Observes nucleotide values, reported in groups of three, in which the position of the third nucleotide is divisible by three. Performed on DNA or RNA.
AA	Amino Acid	Observes amino acid sequences or values. Performed on proteins, or inferred from nucleotide (or codon) results.
ALE	Allele	Identifies the allele (version) of the gene/genetic region of interest. The result is the name for that allele according to the gene's nomenclature committee.

- b. --STRESC holds a "standardized" result. The usual cases where STRESC is different from ORRES are unit conversion and scoring of results (as in questionnaires). PF uses a different convention, where STRESC is the result in a format from an established nomenclature. This standard result often includes information that is parsed out into several SDTM variables, such as those for genetic location (PFGENLOC), observed result (PFORRES), and expected result (PFORREF).
- Nucleotide and amino acid location (position) values in PFGENLOC should follow the same rules as in PFSTRESC.
  - Unless a more appropriate nomenclature exists, the standard format for nucleotide and amino acid results in PFSTRESC follows the rules of HGVS nomenclature.
  - When PFTESTCD is ALE, results in PFSTRESC should follow the nomenclature system specified by the relevant gene's committee.

HGVS is a controversial subject in research genomics!



# If You Can't Name It, You Don't Know It

- Naming is Important!
  - Inevitably, someone is going to set up a database and make a name into a **key**
- Not All Names Are Created Equal
  - normalized?**
  - unique?**
  - structured?**
- External standards may feel like a constraint
  - **But they are a blessing in disguise.**

# Standards vs. Research



**Daniel MacArthur**

@dgmacarthur

+  Follow

Standards are defined by those who have the greatest tolerance for interminably long conference calls about standards.

RETWEETS

88

LIKES

95



6:05 AM - 17 Nov 2015

# What Should You Do?

1. File Formats are not “a detail.”
2. It’s easy to get Naming wrong.
3. Standards Bodies need to include basic researchers as much as possible.

How can we organize our data for both  
**today** and **tomorrow**?

**Please Don't Make  
Your Programmers Cry**