

# Python Based Standardization Tools for ClinicalTrials.Gov



**Jacob Barhak**

**Austin, Texas**

<http://sites.google.com/site/jacobbarhak/>

**COMBINE 2018**

**Boston University**

**October 8-12, 2018**

## **Abstract**

ClinicalTrials.Gov is a government database that stores clinical trials from around the world. This database is growing fast, partially because some requirements of reporting clinical trial results are now supported by U.S. law. Despite the growth of database, the data stored there is still not widely used for modeling and simulation, partially because the database is still a relatively new tool, and partially because the data there is not standardized. Since data is entered by multiple entities into this semi-structured text based database, and since clinical trials have a large variety, the data is not immediately suitable for modeling. Although the National Library of Medicine scrutinizes this data, the scrutiny level is for human comprehensible data, while modeling requires computer comprehension.

For this reason, there is a need to clean data towards tasks such as disease modeling. This work will discuss a set of Python tools that were used to process ClinicalTrials.Gov and aim towards standardization. The tools parse XML data, index the information for easier processing, and then uses Machine Learning and Natural Language Processing to cluster the data, The clustering algorithm is used to assist a human user look at similar information provided by a Graphical User Interface.

These Python tools were used to extract 21,094 units from 30,763 different clinical trails with results. These quantities show a clear need of standardization to be used in future computer modeling efforts.

# Difficulties Making Medical Data Machine Comprehensible

## Data Source

- Electronic Medical Records
- Public datasets and databases:
  - ClinicalTrials.Gov
  - Physionet
  - Many others ...
- Publications
  - Medical journals in electronic form
  - Web sites
  - Printed
- Knowledge held by physicians

## Difficulty

- Restricted and splintered
- New and needs standardization
  - Fast Growing – supported by law
  - Multiple datasets
  - Assorted
- Free Text – not standardized
  - Made for humans + some cost
  - Unorganized non centralized data
  - Assorted and hard to access
- Not Machine Accessible

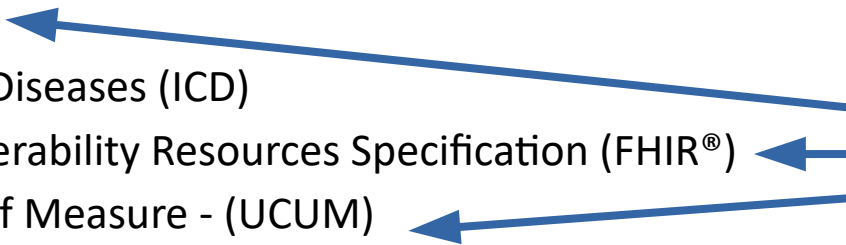
ClinicalTrials.Gov is a good source to start machine comprehension due to: accessibility, size, variety, and potential growth

# Existing Medical Data Specifications/Standards

- Unified Medical Language System (UMLS)
- Snomed CT
- Clinical Data Interchange Standards Consortium (CDISC):
  - Foundational: PRM, SEND, CDASH, SDTM, SDTMIG, SDTMIG-Pgx, ADaM, QRS
  - Data Exchange: CTR-XML, ODM-XML, SDM-XML, Define-XML, DataSet-XML, LAB, RDF
  - Therapeutic Areas – 27 areas including Alzheimer's, Diabetes, Cardiovascular, Influenza, Cancers
  - CDISC Share
  - Semantics
  - Domain information Model
  - Has units classification
- Intentional Classification of Diseases (ICD)
- HL7 Fast Healthcare Interoperability Resources Specification (FHIR®)
- The Unified Code for Units of Measure - (UCUM)

So many specifications!  
Yet, so much data is still not interchangeable  
More work is needed!

Unit  
Standardization



# About ClinicalTrials.gov

- National Institutes of Health (NIH) Project
- Maintained by the National Library of Medicine (NLM)
- Accumulates Clinical Trial Data internationally
- Fast growing database

<b>Date</b>	<b>Number of Trials</b>	<b>Studies with Results</b>
20-Apr-2018	271,510	30,763
3-Nov-2017	258,046	28,785
12-Feb-2017	236,687	24,251
27-Sep-2016	226,460	22,614
7-Apr-2015	187,653	Not collected

- Now many clinical trials are required to register in this data base by law: PUBLIC LAW 110–85—SEPT. 27, 2007 - TITLE VIII—CLINICAL TRIAL DATABASES. Section 801 of the Food and Drug Administration Amendments Act of 2007. Online: <https://www.gpo.gov/fdsys/pkg/PLAW-110publ85/pdf/PLAW-110publ85.pdf#page=82>

# ClinicalTrials.Gov Capabilities

- Has web entry system to collect data from multiple sources
- Entered data is scrutinized for accuracy **at the human level**
- Stores numeric trial results as well as clerical information
- Has a useful **human** web interface
  - Has a sophisticated search capability that handles medical terms, acronyms and filtering
  - Offers expert search for programmers
- Trial data can be exported to XML **for machine readability**
- XML scheme is published

ClinicalTrials.Gov is a great new tool !

Yet, it requires standardization for machine comprehension

# Imported Data Can be Used for Population Generation

Data imported from ClinicalTrials.Gov can be used for generation of synthetic population to **mimic statistics**

- Heterogeneity = generate individuals
- Multiple characteristics per individual
- Allow correlations
- Allow restrictions

Correlated

IndividualID	Male	Age	BP	...
0	0	50	140	...
1	1	45	135	...
2	0	22	120	...
3	1	85	145	...
4	1	14	125	...

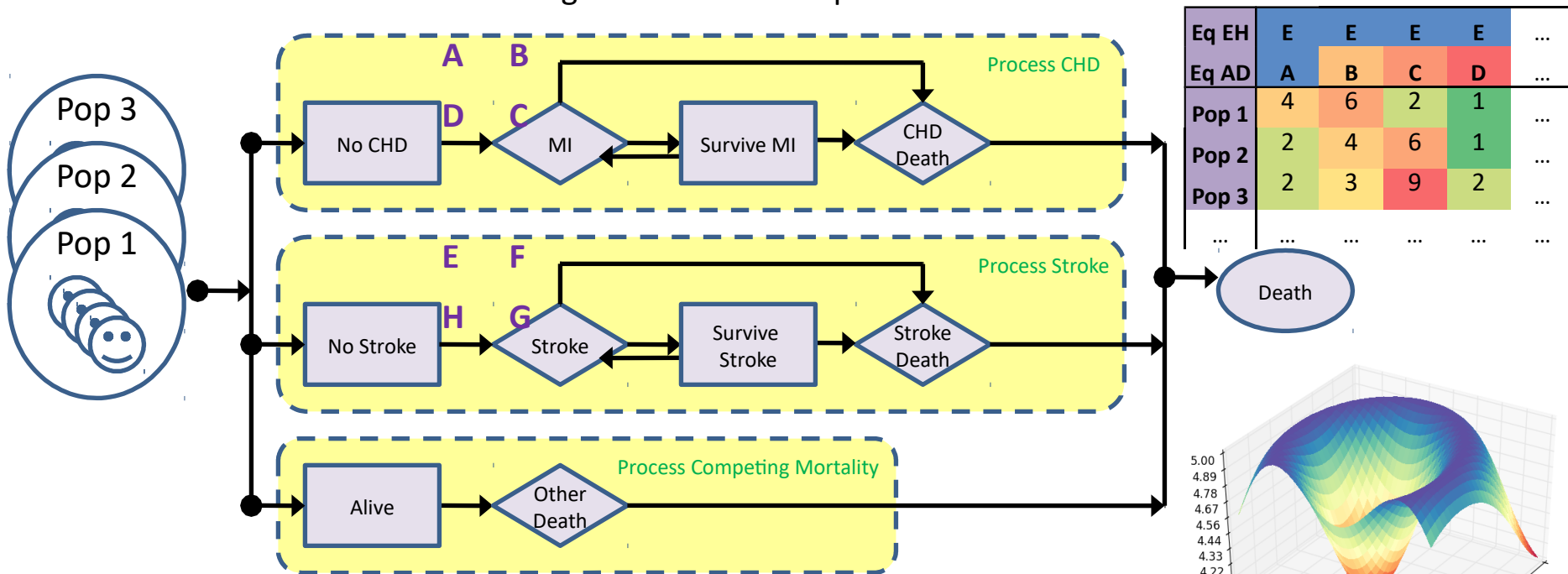
Generation Code / Equations

Restrict Age

# Disease Model Validation with Simulation

## Using Outcomes

- The Reference Model for Disease Progression is an ensemble model
- The model accumulates knowledge: models= assumptions and facts = trial data

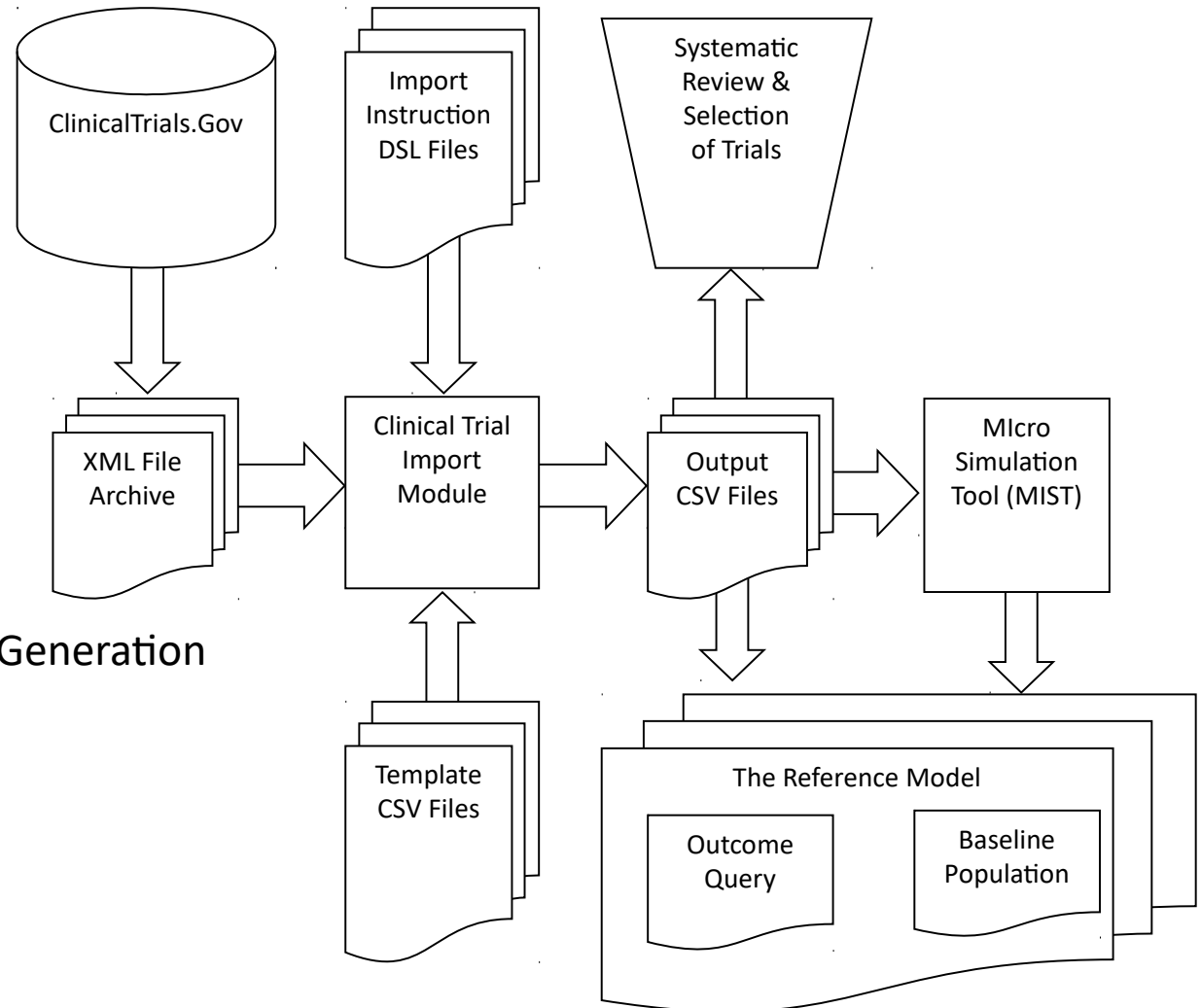


- Model combinations are validated against trial outcomes
- The Reference Model can calculate best model combination

# Importing ClinicalTrials.Gov Data

## Three Import stages:

- Systematic Review
  - Human Selection
- Import Populations
  - Synthetic Population Generation
- Import Outcomes
  - Model Validation





# Issues Encountered During Import

- Systematic review
  - Organizing results in tabular human readable format
  - Extracting meaningful numbers from convoluted free text
- Population generation
  - Code generation from data
  - Name matching
  - Unit conversion – context sensitive
  - Time extraction from free text
- Outcome conversion
  - Scaling numbers to similar reference
  - Cohort matching
  - Calculation of missing outcomes for full trial

The human natural language used in ClinicalTrials.Gov is harder for machines to understand and process!

More importantly, machines cannot interpret numbers without scaling units!

# Python Based Standardization of ClinicalTrials.Gov

- The entire database with results was downloaded on 20-Apr-2018
  - 30,763 XML files of trials were loaded
- A set of python scripts were written with the following purposes:
  - Import XML data and organize it
    - 61 batches were required to handle the big data
  - Index all XML tags and their values
    - 387 different tags = fields were processed
  - Index unit tags with association to title tags
    - 21,094 different units were detected
  - Find similarity between unit names
    - Include CDISC units in similarity matrix
  - Cluster units by similarity
    - 110 clusters were created
  - Create Graphical User Interface for users to manually map units

Several stages required work in smaller batches to fit into memory.

This issue should gain attention in the future due to database growth

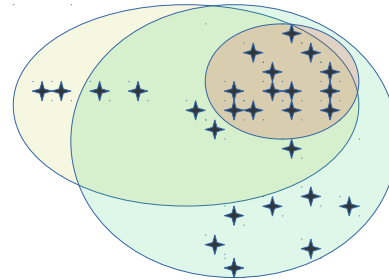
# Natural Language Processing (NLP) Issues

- In addition to ClinicalTrials.Gov units, 6645 CDISC units and their synonyms, including UCUM mapping were added for auxiliary matches
- Unicode characters required special handling
- Synonyms of common words were replaced in pre-processing to help the system with matching, e.g. 'Percent' = '%'
- Distances between text of units was were using two different functions:
  - difflib - calculating charter changes to match words
  - scikit-learn - Term Frequency–Inverse Document Frequency (TFIDF) calculated using 3-6 characters combined with cosine similarity
- Similar units were further ranked:
  - CDISC units were considered more important
  - Common units are ranked first if similar to rare units

Those python libraries produced satisfactory results. An attempt was made with the spaCy library that parses text yet it did not produce better results than the tools used.

# Clustering Units

- To help a user see similar units together, clustering methods were used
- The similarity score between units was used as a metric to cluster similar units
- The MiniBatchKMeans clustering function within scikit-learn was used to reduce memory footprint
- Three clustering passes were made with variations of the similarity score to handle both loose neighbors and very close neighbors



- Clusters in different passes were intersected to capture very close neighbors together
- After clustering, small clusters were merged to larger clusters to improve clustering results by re-merging close neighbors that were split

# Graphic User Interface (GUI)

A user can summon specific cluster #

A user can explore clustering level

A user can save mapping with version #

Suggested synonym

Unit context in trial

#	Used	Unit	Unicode	Mapping	Copy	Possible Synonym	Context	web	Loc
19	8179	1	&#956;g&#183;day/mL	µg·day/mL	Copy	µg·day/mL/mg	tion-time Curve Over the Dosing Interval (AUC&#964;); 'NCT00816400'	Web	Loc
20	4082	3	ng day/mL	ng day/mL	Copy	ng x day/mL	1, 'AUC 0-91 (Area Under Curve)'; 'NCT01143207'	Web	Loc
21	20592	1	ug*H/dL	ug*H/dL	Copy	mg*H/dL	1, 'Cortisol Total Area Under the Curve'; 'NCT01174576'	Web	Loc
22	20591	1	ug*H/L	ug*H/L	Copy	ug*H/L	1, 'Pharmacokinetics: AUC-24'; 'NCT00597493'	Web	Loc
23	2477	5	ug*H/L	ug*H/L	Copy	ug*H/L	salamine Major Metabolite (Ac-5-ASA) at Steady State'; 'NCT01130844'	Web	Loc
24	8108	1	&#181;U hr/mL	µU hr/mL	Copy	ng hr/mL	ge From Baseline in 3-hour Insulin Total AUC at Week 4'; 'NCT02004886'	Web	Loc
25	8125	1	&#181;g times hr/mL	µg times hr/mL	Copy	nanograms times hr/milliliter (ng*hr/mL)	Area Under the Curve (AUC(0 to Infinity)) for Etoricoxib'; 'NCT00945035'	Web	Loc
26	2517	4	&#956;U&#8226;hr/mL	µU·hr/mL	Copy	IU·hr/L	ured by Meal Tolerance Testing (AUC(0-2)) (Final Visit)'; 'NCT01263509'	Web	Loc
27	1769	6	&#956;U&#183;hr/mL	µU·hr/mL	Copy	µg·hr/mL	ured by Meal Tolerance Testing (AUC(0-2)) (Final Visit)'; 'NCT01263496'	Web	Loc
28	6645	2	ng x hr/mL	ng x hr/mL	Copy	ng x h/mL	for AZD6244, 100 mg Administered Orally Twice Daily'; 'NCT00551070'	Web	Loc
29	4531	2	&#181;U* h/mL	µU* h/mL	Copy	IU* h/mL	1, 'Insulin AUECO-5'; 'NCT00558571'	Web	Loc
30	8109	1	&#181;U* h/mL	µU* h/mL	Copy	IU* h/mL	Hours, Following Oral Glucose Tolerance Test ( OGTT )'; 'NCT00605475'	Web	Loc
31	8173	1	&#956;g * hr/mL	µg * hr/mL	Copy	pg * hr/mL	rea Under the Curve (AUC(0 to Infinity)) for Metformin'; 'NCT00929201'	Web	Loc
32	7270	2	pg * hr/mL	pg * hr/mL	Copy	µg * hr/mL	1, 'AUC (0-last) of Active Metabolite GSK614917'; 'NCT00464568'	Web	Loc
33	16971	1	nh.h/mL	nh.h/mL	Copy	ng.h/mL	-21 Days (PK Parameter) Measured for the ASPE Group'; 'NCT01044056'	Web	Loc
34	2508	4	%D.h/mL	%D.h/mL	Copy	%D/mL	1, 'Area Under the Curve (AUC) at 0 to 120 Hours'; 'NCT00315731'	Web	Loc
35	4523	2	&#181;L.h/mL	µL.h/mL	Copy	nmol.h/mL	rea Under Curve (AUC) of Cremophor on Cycle 1 Day 1'; 'NCT01489826'	Web	Loc
36	3230	4	ug.h/L	ug.h/L	Copy	ug.	ve From Time 0 to Infinity of Fospropofol (AUC(0-inf)); 'NCT01260142'	Web	Loc
37	2522	4	&#956;g.h/mL	µg.h/mL	Copy	ng.h/mL	ig Dose AZD9773 Serum Total Fabs (Cohorts 3, 4 and 5)'; 'NCT00615017'	Web	Loc
38	7278	2	pg.h/ml	pg.h/ml	Copy	pg.h/mL	1 Concentration-time Curve From Time Zero to Infinity'; 'NCT02948582'	Web	Loc
39	19763	1	pg.h/mL	pg.h/mL	Copy	pg.h/ml	1, 'AUC0-t'; 'NCT02281448'	Web	Loc
40	14944	1	hr.&#956;g/mL	hr.µg/mL	Copy	µg/mL	Infinity Pharmacokinetic Parameters for a 50 mg Dose'; 'NCT02158936'	Web	Loc
41	16449	1	min&#215;&#956;g/mL	min×µg/mL	Copy	µg/mL	'Area Under the Curve of Plasma Naproxen From 0 to t'; 'NCT00966641'	Web	Loc
42	8166	1	&#956;&#65279;g/mL	µg/mL	Copy	CDISC-UCUM:C64566;g/mL	1, 'Tobramycin Serum Concentration'; 'NCT00918957'	Web	Loc
43	357	37	&#956;g/L	µg/L	Copy	CDISC-UCUM:C42576;g/L	1, 'Baseline PSA'; 'NCT01605227'	Web	Loc
44	2901	4	hr*&#956;g/L	hr*µg/L	Copy	hr*µg/mL	e Extrapolated to Infinity (AUCINFobs) for Etelcalcetide'; 'NCT01134562'	Web	Loc
45	8201		&#956;g/mL x h	µg/mL x h	Copy	µg/mL x h	6 Hour (AUC(0-6)) of Amelastatin B-18 (AmB-18)'; 'NCT01333436'	Web	Loc

A user can map the unit using free text  
This allows mapping of similar units

Shows number of unit occurrences

Show Unicode Escape

Shows the unit as user sees it

# The GUI Assists Human User

For a human, those units are the same while a machine requires exact match

#	Used	Unit	Unicode	Copy	Possible Synonym	Context	web	Loc
174	16160	1	micromole/min/100mg wet weight	micromole/min/100mg wet weight	Copy	microgram/min/100ml	1, 'Citrate Synthase Activity', 'NCT00858845'	Web Loc
175	7481	2	umol/min/100g	umol/min/100g	Copy	ml/min/100g	1, 'Myocardial Fatty Acid Uptake Rate', 'NCT02563834'	Web Loc
176	4098	3	ng/minute	ng/minute	Copy	CDISC:C85749:ng/min	Participants Who Had Normal Glucose Tolerance (NGT)', 'NCT00871871'	Web Loc
177	1949	6	mEq/min	mEq/min	Copy	mEq/minute	Nesiritide Treatment Compared to Placebo Treatment', 'NCT00387621'	Web Loc
178	15510	1	ml*min	ml*min	Copy	ng/mL*min	Under the Effect Curve for Gallbladder Volume (AUEC VL)', 'NCT02496221'	Web Loc
179	4556	2	&#956;L/min	µL/min	Copy	CDISC:C67388:L/min	1, 'Aqueous Flow', 'NCT00572936'	Web Loc
180	10154	1	L/m&#178;	L/m <sup>2</sup>	Copy	kg/m <sup>2</sup>	1, 'Apparent Volume of Distribution', 'NCT00560794'	Web Loc
181	1244	9	L/minute	L/minute	Copy	ml/minute	Respiratory Flow (PEF) Over the 12-Week Treatment Period', 'NCT02040779'	Web Loc
182	1084	11	ml/minute	ml/minute	Copy	ml/minute	1, 'estimated creatinine clearance (eCLCr)', 'NCT02569086'	Web Loc
183	8113	1	&#181;L/min	µL/min	Copy	CDISC:C67388:L/min	Baseline in Retinal Blood Flow After Alikiren or Irbesartan', 'NCT00660309'	Web Loc
184	13495	1	VO2 mL/kg/min	VO2 mL/kg/min	Copy	CDISC-UCUM:C73760:mL/kg/min	1, 'Peak oxygen uptake', 'NCT01009099'	Web Loc
185	9038	1	Change in ml/kg/min	Change in ml/kg/min	Copy	Change in mg/kg/min	1, 'Change in Cardiovascular Fitness', 'NCT02609672'	Web Loc
186	6334	2	mL/Kg	mL/Kg	Copy	CDISC:C67411:mL/kg	1, 'Expiratory Tidal Volume', 'NCT01389882'	Web Loc
187	215	75	mL/kg	mL/kg	Copy	CDISC:C67411:mL/kg	4, 'Volume of Distribution at Steady State (Vss)', 'NCT00496262'	Web Loc
188	6371	2	mcg/min	mcg/min	Copy	CDISC-Synonym:C71211:mcg/min	4, 'Volume of Distribution at Steady State (Vss)', 'NCT00496262'	Web Loc
189	8213	1	&#956;mol/kgFFM/min	µmol/kgFFM/min	Copy	umol/kgFFM/min	2, 'Fluid Balance', 'NCT01612676'	Web Loc
190	4051	3	mmHg*min	mmHg*min	Copy	CDISC:C150900:mmHg*min/L	2, 'Volume at Steady State (Vss: ...hromogenic Assay)', 'NCT01181128'	Web Loc
191	16534	1	ml/min/m^2 BSA	ml/min/m <sup>2</sup> BSA	Copy	ml/min/m <sup>2</sup>	1, 'Apparent Volume of Distribution (Vz)', 'NCT00926263'	Web Loc
192	15235	1	lbs/inch^2	lbs/inch <sup>2</sup>	Copy	lb/in <sup>2</sup>	1, 'Apparent Volume of Distribution...t Day 1 and Day 28', 'NCT00937326'	Web Loc
193	17476	1	pM*min/ml	pM*min/ml	Copy	pM*min/ml	1, 'Change in Extravascular Lung Water (EVLW)', 'NCT00796419'	Web Loc
194	6533	2	min&#8729;pg/ml	min-pg/ml	Copy	CDISC-UCUM:C67327:pg/ml	1, 'Evaluation of Tidal Volume Bene...th Training Alone...', 'NCT00976352'	Web Loc
195	16537	1	ml/min/month	ml/min/month	Copy	ml/min/m <sup>2</sup>	1, 'Extravascular Lung Water Index', 'NCT01675453'	Web Loc
196	4998	2	IU/kg/month	IU/kg/month	Copy	µg/kg/month	1, 'Geometric Mean Apparent Vol...learance Method', 'NCT01393964'	Web Loc
197	8200	1	&#956;g/kg/month	µg/kg/month	Copy	IU/kg/month	1, 'Mean Volume of Distribution (Vz) Post-Single Dose', 'NCT00441337'	Web Loc
198	6323	2	mGFR mL/min	mGFR mL/min	Copy	mGFR mL/min	1, 'Circulating Gastric Inhibitory Polypeptide (GIP) Levels', 'NCT01520454'	Web Loc
199	15491	1	mGFR ml/min	mGFR ml/min	Copy	mGFR mL/min	1, 'Gastric and Intestinal Cell Function in Patients With Different BMI', 'NCT01610154'	Web Loc
200	7480	2	umol/kgFFM/min	umol/kgFFM/min	Copy	umol/kgFFM/min	1, 'Estimated Glomerular Filtration Rate (eGFR)', 'NCT0068107'	Web Loc

CDISC specifications and synonyms are shown to the user

User can view multiple contexts of the same unit and focus on those online

A user can go online focusing on the unit in the trial web page

# The GUI Helps Save Time by Providing Suggestions

Dialog

Cluster: 19 Cluster Pass: -99 Save without exit OK Cancel

#	Used	Unit	Unicode	Mapping	Copy	Possible Synonym	Context	web	Loc
1	15863	1	mg/dl (MAGE)	mg/dl (MAGE)	mg/dl	CDISC-UCUM:C67015:mg/dl	1, 'Glycemic Variability', 'NCT01463878'	Web	Loc
2	15839	1	mg/dL of IgG4	mg/dL of IgG4	mg/dL	CDISC-UCUM:C67015:mg/dl	ut From Baseline Until Desensitization Food Challenge', 'NCT01814241'	Web	Loc
3	15866	1	mg/dl blood	mg/dl blood	mg/dl	mg/dl TG	1, 'Blood Glucose Concentration', 'NCT01137773'	Web	Loc
4	8305	1	(mg/dl)^2/hr^day-1	(mg/dl)^2/hr^day-1	(mg/dl)	CDISC-UCUM:C67015:mg/dL	1, 'Glycemic Labilty Index (GLI)', 'NCT00812487'	Web	Loc
5	15779	1	mg&#708;2/dL&#708;2	mg^2/dL^2	mg/dL	CDISC:C67015:mg/dL	: (Ca&#215;P) Product Levels (mg&#708;2/dL&#708;2)', 'NCT01368042'	Web	Loc
6	3987	3	mg^2/dL^2	mg^2/dL^2	Mg/dL	mg/dL	n (Albumin-adjusted)-Phosphorus Product at Week 22', 'NCT00853242'	Web	Loc
7	4543	2	&#181;ml/ml^day	um/ml^day	mg/dL *h	mg/dL blood	n (Albumin-adjusted)-Phosphorus Product at Week 22', 'NCT00853242'	Web	Loc
8	2361	5	mm^day	mm^day	um/ml	um/ml	dministered on Day 1 of Cycle 1 in Phase I of the Study', 'NCT00517530'	Web	Loc
9	19738	1	percentage^days	percentage^days	mcM^day	mcM^day	General Well-being - Immediate Postoperative Period', 'NCT01355523'	Web	Loc
10	15855	1	mg/day P. Eq.	mg/day P. Eq.	days	days	Area Under The Curve (AUC) on Day 180 and Day 360', 'NCT00162266'	Web	Loc
11	15860	1	mg/day x 5 days	mg/day x 5 days	CDISC:C67399:mg/day	CDISC:C67399:mg/day	Patients Who Received Oral Corticosteroids at Baseline', 'NCT00207740'	Web	Loc
12	5834	2	[ln(mg/dL)]^1.084	[ln(mg/dL)]^1.084	mg twice daily X 5 days	mg twice daily X 5 days	1, 'Maximum Tolerated Dose (Phase I)', 'NCT00404248'	Web	Loc
13	15313	1	ln(mg/L)	ln(mg/L)	ln(mg/L)	ln(mg/L)	1, 'High Blood Glucose Index (HBGI)', 'NCT02137512'	Web	Loc
14	15882	1	mg/kg Q2W	mg/kg Q2W	ln(mg/mL)	ln(mg/mL)	1, 'Plasma C-reactive Protein (CRP)', 'NCT00833898'	Web	Loc
15	15877	1	mg/kg Alefacept	mg/kg Alefacept	Copy	CDISC-UCUM:C67401:mg/kg	1, 'Part 1: PRD of Bevacizumab for Part 2', 'NCT00925769'	Web	Loc
16	8199	1	&#956;g/kg s	ug/kg s	Copy	mg/kg DAC	1, 'Maximum Tolerated Dose (MTD)', 'NCT00438802'	Web	Loc
17	1173	10	mg Fe/g dw	mg Fe/g dw	Copy	ug/kg	1, 'Maximum Tolerated Dose (MTD) of Sylatron', 'NCT01496807'	Web	Loc
18	8381	1	.&#181;g/kg	ug/kg	Copy	mg Fe/g	aseline (BL) to End of Extension Study, by LIC Category', 'NCT00171301'	Web	Loc
19	15902	1	mg/kg2	mg/kg2	Copy	ug/kg	1, 'INTRA OPERATIVE PHENYLEPHRINE USED', 'NCT01069562'	Web	Loc
20	4561	2	&#956;g/kg	ug/kg	Copy	CDISC-UCUM:C67401:mg/kg	1, 'Body mass index', 'NCT02129725'	Web	Loc
21	15893	1	mg/kg, Units/kg	mg/kg, Units/kg	Copy	CDISC-UCUM:C69104:g/kg	1, 'Darbepoetin Alfa Weight-Adjusted Dose Over Time', 'NCT00436748'	Web	Loc
22	2969	4	mg/kg, units/kg	mg/kg, units/kg	Copy	mg/kg, units/kg	Short-Acting Insulin Dose (Adjusted for Body Weight)', 'NCT00138671'	Web	Loc
23	2514	4	&#181;g/kg	ug/kg	Copy	mg/kg, unit/kg	ly Short-Acting Insulin Dose Adjusted for Body Weight', 'NCT00139659'	Web	Loc
24	194	85	mg/kg	mg/kg	Copy	.ug/kg	1, 'Cumulative Dose of Norepinephrine', 'NCT01612676'	Web	Loc
25	2962	4	mg&#178;/dL&#178;	mg^2/dL^2	Copy	CDISC-UCUM:C67401:mg/kg	0 (Induction period) or week 16 (Maintenance period)', 'NCT01393405'	Web	Loc
26	15777	1	mg&#12539;h / dL	mg^h / dL	Copy	CDISC-UCUM:C67015:mg/dL	4, 'Corrected Calcium Phosphorus Product (Ca x P)', 'NCT01785875'	Web	Loc
27	6310	2	mg/dL	mg/dL	Copy	mg^h / dL	the Curve (AUC) 0 to 2h (Breakfast, Lunch and Dinner)', 'NCT01072331'	Web	Loc
					Copy	CDISC-UCUM:C64572:mg/dL	ate Postprandial LFA at Week 12 in the 0-40 mg Cohort', 'NCT00650502'	Web	Loc

User can copy selected suggested unit synonym to mapping field

User can explore similar suggested units and pick more suitable mapping

# Very Common Units and Terms

**Units appearing more than 1000 times:** Participants , participants , years , units on a scale , percentage of participants , Percentage of participants , months , Units on a scale , Years , mg/dL , ng/mL , days , Subjects , mmHg , scores on a scale , Percentage of Participants , Scores on a scale , percent change , hours , Months , Days , mm , kg , mmol/L , subjects , minutes , pg/mL , score on a scale , ratio , kg/m^2

#	Used	Unit	Unicode	Mapping	Copy	Possible Synonym	Context	web	Loc
88	5773	2	U/mg protein	U/mg protein	Copy	mg protein	1, 'Change in Glutathione Peroxidase Activity', 'NCT02331446'	Web	Loc
89	3980	3	mg protein	mg protein	Copy	U/mg protein	1, 'Wheat OFC Dose at First Symptom', 'NCT01980992'	Web	Loc
90	19392	1	percentage of recording failure	percentage of recording failure	Copy	% of recordings	1, 'Quality of Polysomnographic Recordings', 'NCT01471626'	Web	Loc
91	19101	1	percentage of no PFS failure	percentage of no PFS failure	Copy	% of no PSA failure	ing 3 Years Treatment Between Degarelix and Goserelin', 'NCT01242748'	Web	Loc
92	19102	1	percentage of no PSA failure	percentage of no PSA failure	Copy	% of no PFS failure	ing 3 Years Treatment Between Degarelix and Goserelin', 'NCT01242748'	Web	Loc
93	2030	6	percentage of kilocalorie intake	percentage of kilocalorie intake	Copy	% of kilocalories	ervey From Carbohydrate at the Meal (Exercise Session)', 'NCT00943436'	Web	Loc
94	7225	2	percentage of time spent at goal	percentage of time spent at goal	Copy	% of time spent in range	1, 'Primary Effectiveness Outcome - A1c', 'NCT01586897'	Web	Loc
95	19575	1	percentage of time in 12 h after drug	percentage of time in 12 h after drug	Copy	% of time in bed	r Oximetry Saturations Under 90% in 2-6 Month Infants', 'NCT01260883'	Web	Loc
96	19584	1	percentage of time spent < 65 mg/dl	percentage of time spent < 65 mg/dl	Copy	% of time spent > 180 mg/dl	ercentage of Time Spent at Glycemic Levels <65 mg/dl', 'NCT01341067'	Web	Loc
97	19585	1	percentage of time spent > 180 mg/dl	percentage of time spent > 180 mg/dl	Copy	% of time spent < 65 mg/dl	ercentage of Time Spent at Glycemic Levels > 180 mg/dl', 'NCT01341067'	Web	Loc
98	18084	1	percent of time spent in hypoglycemia	percent of time spent in hypoglycemia	Copy	% of time spent in hospital	1, 'Percentage of Time Spent in Hypoglycemia', 'NCT00606034'	Web	Loc
99	19655	1	percentage of tubes retained	percentage of tubes retained	Copy	% of tubes	1, 'Tube Retention', 'NCT01202578'	Web	Loc
100	7392	2	retained tubes	retained tubes	Copy	% of tubes retained	1, 'Tube Retention', 'NCT00939796'	Web	Loc
101	18586	1	percentage of accurate diagnoses	percentage of accurate diagnoses	Copy	% of accurate taps	age of Participants With Accurate Diagnoses of Cancer', 'NCT01227382'	Web	Loc
102	19537	1	percentage of surgery minutes	percentage of surgery minutes	Copy	% of surgeries	ve Case Time With Systolic Blood Pressure <95 mmHg', 'NCT01208402'	Web	Loc
103	19137	1	percentage of ostia treated	percentage of ostia treated	Copy	% of ostia patent	1, 'Patency of Treated Area', 'NCT01623050'	Web	Loc
104	19350	1	percentage of preoperative values	percentage of preoperative values	Copy	% of change negative values	1, 'Change From Baseline of Spirometric Values', 'NCT01249872'	Web	Loc
105	18695	1	percentage of breast density	percentage of breast density	Copy	% of breasts	1, 'Change in Percent Density', 'NCT00859651'	Web	Loc
106	1719	7	percentage of opiate negative	percentage of opiate negative	Copy	% of true negatives	1, 'Percent Opiate Negative', 'NCT00249470'	Web	Loc
107	19453	1	percentage of specimens ESA negative	percentage of specimens ESA negative	Copy	% of specimen	nor Specimens Nonreactive by ABBOTT PRISM Chagas', 'NCT01662362'	Web	Loc
108	19525	1	percentage of submitted negative samples	percentage of submitted negative samples	Copy	% of true negative sample	ge Samples Submitted Negative for Crack Cocaine Use', 'NCT01815645'	Web	Loc
109	17379	1	number of true negatives	number of true negatives	Copy	% of true negatives	1, 'Specificity', 'NCT00620373'	Web	Loc
110	18872	1	percentage of efficient sleep	percentage of efficient sleep	Copy	% of efficiency	1, 'Sleep Efficiency (SE)', 'NCT00374556'	Web	Loc
111	7097	2	percentage of fat in liver	percentage of fat in liver	Copy	% of fat intake	r Fat by Magnetic Resonance and Spectroscopy (MRS).', 'NCT00994682'	Web	Loc
112	5559	2	Percentage of technical events	Percentage of technical events	Copy	number of technical successes	idence of Technical Events Experienced by Participants', 'NCT01956032'	Web	Loc
113	19549	1	percentage of taking adherence	percentage of taking adherence	Copy	% of adherence	requency of Medication-taking (% Taking Adherence)', 'NCT02553512'	Web	Loc
114	7017	2	percentage of adhered platelets	percentage of adhered platelets	Copy	% coated platelets	1, 'Platelet Adhesion 2 Hours', 'NCT02140372'	Web	Loc

A large part of the units currently used are very similar and contain outcome text

Common terms are considered more important due to majority reporting

Many numbers are proportions and in that sense unit-less, yet represented as percentage So the unit text represents scaling



# Standardization is Important

#	Used	Unit	Unicode	Mapping	Copy	Possible Synonym	Context	web	Loc
36	18015	1	percent of hemoglobin glycosylated	percent of hemoglobin glycosylated	Copy	% of hemoglobin glycosylated	1, 'HbA1c', 'NCT02157298'	Web	Loc
37	18966	1	percentage of hemoglobin glycosylated	percentage of hemoglobin glycosylated	Copy	% of hemoglobin glycosylated	1, 'Adjusted Mean Change in HbA1c Levels', 'NCT02157298'	Web	Loc
38	446	28	Percent (%) glycosylated haemoglobin	Percent (%) glycosylated haemoglobin	Copy	% of glycosylated haemoglobin	4, 'Glycosylated Haemoglobin (HbA1c)', 'NCT00474045'	Web	Loc
39	11767	1	Percent of glycosylated hemoglobin	Percent of glycosylated hemoglobin	Copy	% of glycosylated hemoglobin	1, 'Hemoglobin A1c (A1C)', 'NCT00885352'	Web	Loc
40	18943	1	percentage of glycosylated hemoglobin	percentage of glycosylated hemoglobin	Copy	% of glycosylated hemoglobin	1, 'Change in Hemoglobin A-1C From Baseline', 'NCT00254501'	Web	Loc
41	18944	1	percentage of glycosylated hemoglobin	percentage of glycosylated hemoglobin	Copy	% of glycosylated hemoglobin	1, 'HbA1c', 'NCT00351819'	Web	Loc
42	1413	8	Percent of glycosylated hemoglobin (A1C)	Percent of glycosylated hemoglobin (A1C)	Copy	% of glycosylated hemoglobin	2, 'Hemoglobin A1C (A1C)', 'NCT00541775'	Web	Loc
43	18945	1	percentage of glycosylated hemoglobin	percentage of glycosylated hemoglobin	Copy	% of glycosylated hemoglobin	1, 'HbA1c at Week 14', 'NCT00351819'	Web	Loc
44	141	139	percentage of glycosylated haemoglobin	percentage of glycosylated haemoglobin	Copy	% of glycosylated haemoglobin	40, 'Glycosylated haemoglobin (HbA1c)', 'NCT00612040'	Web	Loc
45	738	16	Percentage of glycosylated haemoglobin	Percentage of glycosylated haemoglobin	Copy	% of glycosylated haemoglobin	2, 'Glycosylated haemoglobin (HbA1c)', 'NCT01850615'	Web	Loc
46	11766	1	Percent of glycosylated haemoglobin	Percent of glycosylated haemoglobin	Copy	% of glycosylated haemoglobin	From Baseline in Hemoglobin A1C (HbA1c) at Week 24', 'NCT01619059'	Web	Loc
47	3145	4	percentage of glycosylated hemoglobins	percentage of glycosylated hemoglobins	Copy	% of glycosylated hemoglobins	moglobin A1c (A1C) at Week 12 in the 0-40 mg Cohort', 'NCT00950599'	Web	Loc
48	3651	3	Percentage of glycosylated hemoglobins	Percentage of glycosylated hemoglobins	Copy	% of glycosylated hemoglobins	anges From Baseline at Week 24 - Open Label Cohort', 'NCT00121641'	Web	Loc
49	5482	2	Percentage of Glycosylated Haemoglobin	Percentage of Glycosylated Haemoglobin	Copy	% of glycosylated haemoglobin	1, 'Screening HbA1c', 'NCT01519466'	Web	Loc
50	174	102	percentage of Glycosylated Hemoglobin	percentage of Glycosylated Hemoglobin	Copy	% of Glycosylated Hemoglobin	From Baseline in Glycosylated Hemoglobin (Week 12)', 'NCT00286442'	Web	Loc
51	3635	3	Percentage of Glycosylated Hemoglobin	Percentage of Glycosylated Hemoglobin	Copy	% of Glycosylated Hemoglobin	1, 'Glycosylated Hemoglobin (HbA1c)', 'NCT01533428'	Web	Loc
52	5424	2	Percent of Glycosylated Hemoglobin	Percent of Glycosylated Hemoglobin	Copy	% of Glycosylated Hemoglobin	ange From Baseline in HbA1c at Week 12 and Week 24', 'NCT00666718'	Web	Loc
53	669	18	percent glycosylated hemoglobin	percent glycosylated hemoglobin	Copy	% glycosylated hemoglobin	3, 'Glycosylated Hemoglobin (HbA1c)', 'NCT00279201'	Web	Loc
54	1097	11	percentage glycosylated hemoglobin	percentage glycosylated hemoglobin	Copy	% glycosylated hemoglobin	ge From Baseline in Glycosylated Hemoglobin (HbA1c)', 'NCT01289119'	Web	Loc
55	4470	2	% glycosylated hemoglobin	% glycosylated hemoglobin	Copy	% glycosylated hemoglobin	1, 'Change in Hb A1c From Baseline to 12 Months', 'NCT00493012'	Web	Loc
56	200	82	percentage of glycosylated hemoglobin	percentage of glycosylated hemoglobin	Copy	% of glycosylated hemoglobin	8, 'Glycosylated Hemoglobin (HbA1c)', 'NCT00630825'	Web	Loc
57	845	14	Percentage of glycosylated hemoglobin	Percentage of glycosylated hemoglobin	Copy	% of glycosylated hemoglobin	1, 'Baseline Hemoglobin A1c (HbA1c)', 'NCT02302716'	Web	Loc
58	903	13	Percent of glycosylated hemoglobin	Percent of glycosylated hemoglobin	Copy	% of glycosylated hemoglobin	3, 'Hemoglobin A1c (HbA1c)', 'NCT00666718'	Web	Loc
59	1023	12	percent of glycosylated hemoglobin	percent of glycosylated hemoglobin	Copy	% of glycosylated hemoglobin	1, 'Baseline Glycosylated Hemoglobin', 'NCT00420095'	Web	Loc
60	4475	2	% of glycosylated hemoglobin	% of glycosylated hemoglobin	Copy	% of glycosylated hemoglobin	1, 'HbA1c', 'NCT01182493'	Web	Loc
61	18023	1	percent of length at end of filling	percent of length at end of filling	Copy	% of bone filling	3 Systolic Maximum Shortening (LVES Max Shortening)', 'NCT01052272'	Web	Loc
62	18860	1	percentage of each WBC type in WBC count	percentage of each WBC type in WBC count	Copy	% of WBC count	ubocytes, Monocytes, and Total Neutrophils at Week 8', 'NCT00742508'	Web	Loc

So many different ways are used to write the unit of HbA1c

The system suggestion for standardization is reasonable yet the user has to choose what is best

# Standardization Reveals Mistakes

#	Used	Unit	Unicode	Mapping	Copy	Possible Synonym	Context	web	Loc
56	12676	1	R&#178;	R²	Copy	CDISC-UCUM:C70575:R	dose Value in Lactate Versus Metformin Concentration', 'NCT01658514'	Web	Loc
57	3466	3	KG	KG	Copy	CDISC-UCUM:C42537:K	1, 'Weight (kg)', 'NCT01626118'	Web	Loc
58	9377	1	ED	ED	Copy	AEDs	rpes Simplex Virus (Anti-HSV) Neutralizing Antibodies', 'NCT00057330'	Web	Loc
59	11415	1	PU	PU	Copy	CDISC:C67264:PFU	Development of Mesenteric Traction Syndrome (MTS)', 'NCT02951273'	Web	Loc
60	2103	5	&#186;C	°C	Copy	CDISC:C42559:C	1, 'Baseline Thermal Pain Threshold', 'NCT02320838'	Web	Loc
61	538	22	Hz	Hz	Copy	CDISC-UCUM:C42545:Hz	1, 'B1-Fmean', 'NCT01955083'	Web	Loc
62	2806	4	U	U	Copy	CDISC:C44278:U	inge in Insulin Glargine Dose From Baseline to Week 26', 'NCT01768559'	Web	Loc
63	1112	10	&#176;C	°C	Copy	CDISC:C42559:C	1, 'Body temperature', 'NCT02048072'	Web	Loc
64	823	14	&#181;V	µV	Copy	CDISC-UCUM:C42551:V	ctinogram (ERG) at Month 12 as Compared to Baseline.', 'NCT01838655'	Web	Loc
65	486	25	&#181;M	µM	Copy	µM-h	kinetic Parameters (Maximum Plasma Concentration)', 'NCT00439218'	Web	Loc
66	2520	4	&#956;V	µV	Copy	CDISC-UCUM:C42551:V	oked Potentials at the Midline Occipital Electrode (Oz)', 'NCT02079844'	Web	Loc
67	8105	1	&#176;F	°F	Copy	CDISC:C44277:F	Mean Change From Baseline in Temperature (&#176;F)', 'NCT00711009'	Web	Loc
68	8523	1	47	47	Copy	mm^4	1, 'Race (NIH/OMB)', 'NCT01494298'	Web	Loc
69	4529	2	&#181;S	µS	Copy	µSv	'Smoking Cues, Measured Using Script Driven Imagery', 'NCT00916721'	Web	Loc
70	2161	5	Gy	Gy	Copy	CDISC-UCUM:C18063:Gy	1, 'Tongue dose', 'NCT01576939'	Web	Loc
71	21073	1	y	y	Copy	CDISC:C48553:yd	1, 'Age', 'NCT01801280'	Web	Loc
72	9754	1	H	H	Copy	CDISC-UCUM:C42558:H	1, 'Change in the Shannon Diversity Index', 'NCT02465463'	Web	Loc
73	9914	1	IQ	IQ	Copy	ISQ	1, 'Intelligence Quotient (IQ)', 'NCT00588731'	Web	Loc
74	8527	1	68	68	Copy	16	1, 'Age', 'NCT01841021'	Web	Loc
75	4554	2	&#8451;	°C	Copy	CDISC:C42549:Watt	ollowed by 30 Min After Brachial Plexus Block(Phase 1)', 'NCT02139982'	Web	Loc
76	4437	2	\$	\$	Copy	US\$	1, 'Gross Weekly Purchasing of Fruits and Vegetables', 'NCT01509664'	Web	Loc
77	8165	1	&#937;	Ω	Copy	CDISC:C42549:Watt	1, 'Effectiveness 3', 'NCT01688843'	Web	Loc
78	8079	1	%B	%B	Copy	CDISC-UCUM:C25613:%	'Homeostatic Model Assessment (HOMA)-%B by Visit', 'NCT00263328'	Web	Loc
79	13538	1	W	W	Copy	CDISC-UCUM:C42556:Wb	1, 'Muscle Strength', 'NCT00929500'	Web	Loc
80	4720	2	C	C	Copy	CDISC:C42559:C	hange From Baseline Temperature on Day 4 and Day 7', 'NCT02074358'	Web	Loc
81	20	1678	Days	Days	Copy	dy	36, 'Age', 'NCT00129129'	Web	Loc
82	321	44	Day	Day	Copy	dy	3, 'Terminal Phase Half-Life (T1/2)', 'NCT00304654'	Web	Loc

These are mistakes in data entry that were not caught during screening

A standardization effort will reveal those errors and improve quality

# Future work

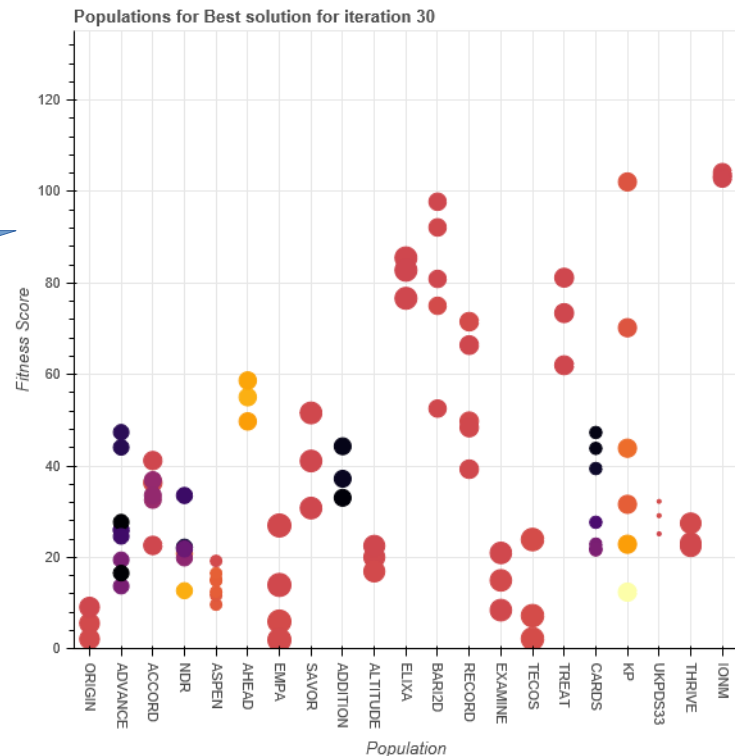
- **In the Foreseeable Future:**

- Standardizing units from ClinicalTrials.Gov into CDISC and UMLS and possible collaboration with the Simulation Interoperability Standards Organization - SISO
- Parallel efforts expanding SBML to represent disease models
- The Reference Model is continuously absorbing new knowledge in the form of models and data to better map our cumulative computational understanding

A map that shows our computational gap of understanding of clinical trials is already possible! Much due to ClinicalTrials.Gov data availability

- **In the Much Farther Future:**

- Computer comprehension of clinical data.
- Allow applications such as:
  - Personal computerized medical adviser
  - Preventative medical predictor



# Acknowledgments

- Thanks to CDISC consortium help
- Thanks to NIH persons who helped and specifically to:
  - Nick Ide from NLM ClinicalTrials.Gov team on advice
  - Erin E Muhlbradt from NCI for advice on CDISC unit data
- Thanks to the SBML community and DDMoRe consortium, specifically to:
  - Chris Myers, University of Utah, USA
  - Leandro Watanabe, University of Utah, USA
  - Lucian Smith, Caltech, USA
  - Jacek Swat, Simcyp Ltd. (Certara), UK

# References

- **The Reference Model and ClinicalTrials.Gov**
  - J. Barhak, The Reference Model Models ClinicalTrials.Gov. SummerSim 2017 July 9-12, Bellevue, WA. Paper: <https://doi.org/10.22360/SummerSim>
  - J. Barhak, The Reference Model Visualizes Gaps in Computational Understanding of Clinical Trials, 2018 IMAG Futures Meeting March 21-22, 2018 @ NIH, Bethesda, MD. [http://sites.google.com/site/jacobbarhak/home/Poster\\_IMAG\\_MSM2018\\_Map\\_Upload\\_2018\\_03\\_17.pdf](http://sites.google.com/site/jacobbarhak/home/Poster_IMAG_MSM2018_Map_Upload_2018_03_17.pdf)
  - J. Barhak, The Reference Model: A Decade of Healthcare Predictive Analytics with Python, PyTexas 2017, Nov 18-19, 2017, Galvanize, Austin TX. Presentation: [http://sites.google.com/site/jacobbarhak/home/PyTexas2017\\_Upload\\_2017\\_11\\_18.pptx](http://sites.google.com/site/jacobbarhak/home/PyTexas2017_Upload_2017_11_18.pptx) Video: [https://youtu.be/Pj\\_N4izLmsI](https://youtu.be/Pj_N4izLmsI)
- **Related Future Standardization Efforts for Models and Data**
  - Jacob Barhak, Chris Myers, Leandro Watanabe, Lucian Smith, Maciek Jacek Swat, Healthcare Data and Models Need Standards. Simulation Interchangeability Standards Organization (SISO) 2018 Fall Innovation Workshop. 9-14 Sep 2018 Orlando, Florida. Presentation: [http://sites.google.com/site/jacobbarhak/home/SISO\\_SIW\\_2018\\_08\\_14.pptx](http://sites.google.com/site/jacobbarhak/home/SISO_SIW_2018_08_14.pptx)
  - L. Smith, M. J. Swat, J.Barhak. Sharing Formats for Disease Models. SummerSim 2016 24-27 July, Montreal, CA. Paper: <https://doi.org/10.22360/SummerSim.2016.SCSC.010>
  - L. Watanabe, J. Barhak, C. Myers, Towards Reproducible Disease Models using the Systems Biology Markup Language. Simulation 2018. <http://dx.doi.org/10.1177/0037549718793214> Accompanying Source code: <https://github.com/Jacob-Barhak/DiseaseModelsSBML>