



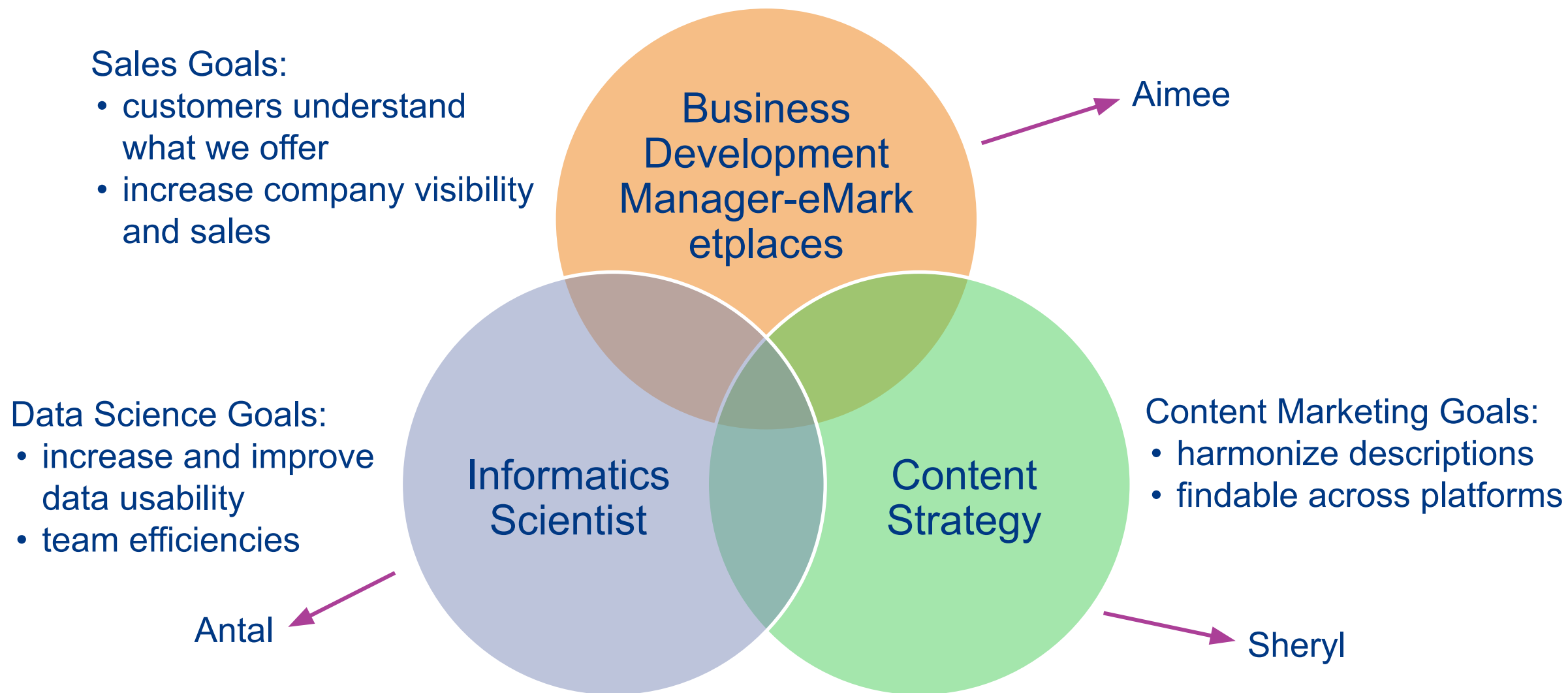
# Describing the Need: An Ontology End-user Case Study

Sheryl Denker & Antal Berenyi

Workshop on Synergizing Biomedical Ontologies  
July 16, 2021



A STRONG  
FOUNDATION  
FOR SUCCESSFUL  
DRUG DISCOVERY



# The Challenge of eMarketplaces

- ~1500 provided terms manually curated by 6 employee contributors (Product/Content Managers)
- Tedious effort resulted in a list of 523 terms
- A second, independent pool of 5K terms was provided
- How can we:
  - ✓ reduce manual labor in the curation effort
  - ✓ obtain only highly relevant results
  - ✓ not miss important terms

Leukemia Models	Capillary Electrophoresis	Human Immortalised Cell Lines	Ocular Permeability	Biosensor Assays	Histone Deacetylase (HDAC)
Thrombosis Models	Cell Free DNA Assay	Human iPSC's	Skin Permeability	Bone Resorption Assay	In Vitro Enzyme Activity Assays
Zucker Rat	Digital Polymerase Chain Reaction	Human Adult Stem Cells	Vaginal Permeability	Catecholamine Assay	Kinase Assays
Bone Metastases Models	DNA Pull Down Assay	Human MSC's	Cell Fractionation	Cell Based Assays Consult	Methyl CpG Binding Domain (MeCPG)
Bone Models	DNA-lipid interaction analysis	Human HSC's	Compound Distribution	Cell Membrane Lipidomics	Methyltransferase Assays
Cartilage Repair Models	DNase Footprinting Assay	Human DNA	Drug Distribution Studies	Chorioallantoic Membrane (CAM)	Mutase Assays
In Vitro Bone Models	Electrophoretic Mobility	Human RNA	Pharmacokinetic - A Pharm	Chromium Release Assay	Nitric Oxide Synthase Assays
Osteoarthritis Models	Isothermal Amplification	Human Cell Lysates	Red Blood Cell (RBC) Parameters	Ciliary Motility	Oxidase Assays
Osteoporosis Models	Linear DNA Amplification	Human Cell Products	Conjugate Detection	Co-cultivation Microscopy Assays	Oxygenase Assays
Artery On-A-Chip Platform	Plasmid Retention Assay	Human Whole Blood	Glucuronidation Assays	Co-Culture Angiogenesis Assays	PARP Assays
Atrial Arrhythmias Models	Polyacrylamide Gel Electrophoresis	Human Umbilical Cord Blood	Glutathione Transferase	Coagulation	PDZ Domain Assays
Cardiovascular Metabolic Models	qPCR Analysis	Human Serum	N-Acetyl Transferase	Cytokine Release Assay	Peroxidase Assays
Cardiovascular Models	Recombinant Nucleoside	Human Plasma	Sulfotransferase	Cytokines	Phosphatase Assays
Coronary Artery Disease Models	Reverse Transfection	Human RBC's	CYP Induction	Ecarin Clotting Time (ECT)	Phosphodiesterase Assays
ExeGen® LDLR MiniSwine	Threshold DNA Detection	Human Buffy Coat	CYP Inhibition Assay	Factor X Assay	Phosphorylase Assays
In Vitro Cardiovascular Models	Tissue Microarray	Human PBMC's	CYP Isozyme Mapping	Fibrinogen Assays	Polymerase Assays
In Vitro Hypoxia Models	Whole Genome Amplification	Human Leukopacks	Cytochrome P450s	Glucose Release Assay	Protease Assay
Ischemia Models	Amplicon - HLA Typing	Human Biofluids	Corning® HepatoCells	Glucose Uptake	Racemase Assays
Japanese White Rabbit	Amplicon 16s rRNA Seq	Human synovial fluid	Drug Drug Interactions	Hairless Mouse	RNA Methylation Assays
Myocardial Infarction Models	ChIP-sequencing Profiling	Human semen	Phase I Metabolism	High Content Screening	Synthetase Assays
Ossabaw-ptFH3 Patient Model™	De Novo Genome Sequencing	Human vaginal fluid/swap	Phase II Metabolism	High Throughput Screening	Ubiquitin Ligase Assays
Pulmonary Hypertension Animal Models	Exome Sequencing	Human amniotic fluid	GI Fluid Stability	Histamine Release Assay	UDP-glucuronosyltransferase
Restenosis Models	Metagenomic Sequencing	Human ascites	Metabolic Stability	In Vivo Angiogenesis Assay	Anaplastic lymphoma kinase (ALK)
SHR/NCrl Rat	mtDNA Sequencing	Human Vitreous Humour	Metabolic Stability in Hep	Inflammation Assays	Angiotensin-2 (Ang-2) Assays
Ventricular Tachycardia Models	NGS: Next Generation Sequencing	Human Tears	Metabolic Stability in Micro	Interferon Gamma Release Assay	B-Cell Lymphoma 2 (BCL-2) Assays
White Rabbit	RAD Sequencing	Human Sputum	Plasma Stability Assay	International Normalized Ratio	Carcinoembryonic Antigen (CEA)
Zecardio	Sanger Sequencing	Human gastrointestinal fluid	Isotopic Labeling	Lipid Detection	CD134 (OX 40) Assays
Zucker Rat	Targeted Sequencing/PCR	Human Bone Marrow aspirate	Metabolite Generation	Lipid Disorder Screening	CD137 Assays
Acute Kidney Failure Models	TCR Repertoire Sequencing	Human Faecal Matter	Metabolite Identification	Lipid Kinase Assays	CD20 Assays
Acute Liver Failure Models	Whole Genome Bisulphite	Human Nasal Fluid	Metabolite Quantification	Lipid Peroxidation Assay	CD252 (OX 40L) Assays
Acute Lung Injury Models	Whole Genome Sequencing	Human Broncho-Alveolar Lavage	Metabolite Structural Elucidation	Lipolysis	CD79b Assays
Critical Care Models	ChIP-sequencing Profiling	Human Breast Milk	Metabolite Synthesis	Low Density Lipoprotein (LDL)	Checkpoint Kinase-1 (CHK-1) Assays
Emesis Models	Expression Profiling	Human Bile	NMR Methods	Matrigel Plug Assay	Colony Stimulating Factor 1 Receptor
Liver Disease Animal Models	RNA Profiling	Human Saliva	No Label Methods	Metabolomics	Cytotoxic T-Lymphocyte Antigen
Liver Fibrosis Models	Target Gene Expression	Human Urine	Non Isotopic Labeling	Micromanipulation Microscopy	Extracellular signal Regulated Kinase
Sepsis Models	Whole Genome Profiling	Human cerebral spinal fluid (CSF)	Radiolabeled Tracers	Oxidative Burst	Fibroblast Growth Factor Receptor
Tissue Rejection Models	CRISPR	Human TMA	Reactive Metabolite	Pathway Mapping	Human Double Minute 2 (HDMDM2)
Zucker Rat	Engineered Meganucleotides	Human custom TMA	Reference Standards	Phenotype	Human Epidermal growth factor
Bone Models	Genome Editing	Human Brain	Equilibrium Dialysis	Phenotypic Assays	Human Organoids
Atopic Dermatitis Models	Recombinant AAV Med	Human Tonsils	Protein Aggregation	Platelet Function Testing	Hypoxia-Inducible Factor 1 (HIF1)
Dermatology Models	Transcription Activator	Human Eyes	Protein Binding HPLC	Proteomics	Immunonecrosis Assays
Hair Growth Models	Zinc Finger Nuclease Genotyping	Human Ear	Protein Binding Ultrafiltration	Reactive Species Assays	Indoleamine 2,3-Di-Oxygenase
Hairless Mouse	Allele Specific Oligonucleotide	Human Hair		Serotonin Release Assay	Kallikrein (KLK)-related Serine
Hairless Rat	Amplified Fragment Length	Human Larynx		Sprouting Assay	Killer cell Immunoglobulin-like
In Vitro Skin Models	DNA Hydroxy Methylation	Human Head and Neck Samples		Stem Cell Differentiation Microarrays	Liposome Immunoassay
Lupus Models	Genotyping	Human Lung		Synthetic Lethality Screening	Lymphocyte Activation Gene-1
Pruritus (Itch)	Random Amplified Polymorphic	Human Heart		Thrombin Generation Tests	Lysine Specific Demethylase-1
Psoriasis Tissue Models	Restriction Fragment Length	Human Vessels		Thrombophilia Screening	Mammalian Target of Rapamycin
Rosacea Models	Artificial Chromosomes	Human Breast		Tube Formation Assay	Mesothelin Assays

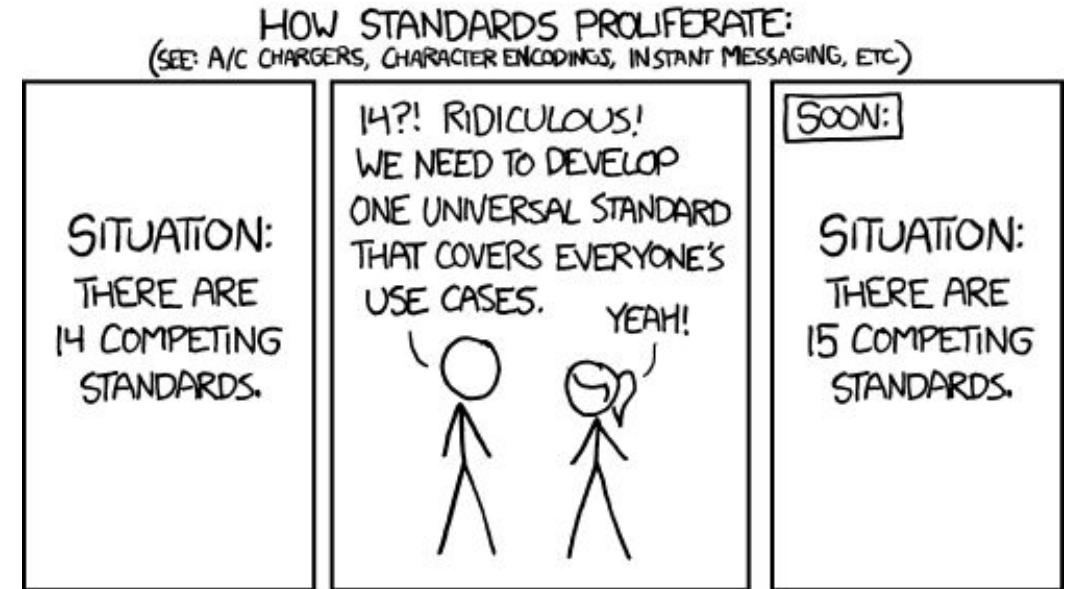


# Use An Ontology! But Which One?

~900 ontologies related to biomedical sciences

Search repository <https://bioportal.bioontology.org>

- cell-based assays
- phenotypic assays



<https://xkcd.com/927/>

20-30 results, including BioAssay Ontology, SNOMEDCT, GO-PLUS

Where do we go from here?

# Let's Use A Mapping Tool!

**OXO**  
ONTOLOGY MAPPING

Home | Documentation | About

### Mapping results

Select a term to see more information. The evidence column tells you how many times we have seen this mapping and the distance is a how many steps need to go to find this mapping. Distance 1 is a direct mapping, the greater the distance the less likely it is that a mapping holds true. Max distance is 10.

Could not map: 521

Input	Mapped id	
Wikipedia:Cosmetics	CHEBI:64857 (cosmetic)	CHEBI
Wikipedia:Inflammation	GO:0006954 (inflammatory response)	GO
3D Cell Culture	No mapping	-
3D Spheroids	No mapping	-

- OxO finds cross-references between ontologies, **vocabularies** and coding standards
- Only 2 out of 523 terms were found: inflammatory response and cosmetic
  - Chemical Entities of Biological Interest (CHEBI)
  - Gene Ontology (GO)

Finding and using the correct ontology is challenging for end users.  
Our ultimate goal is to contribute assay data to an ontology that we find helpful.

## 1. Biology and Users

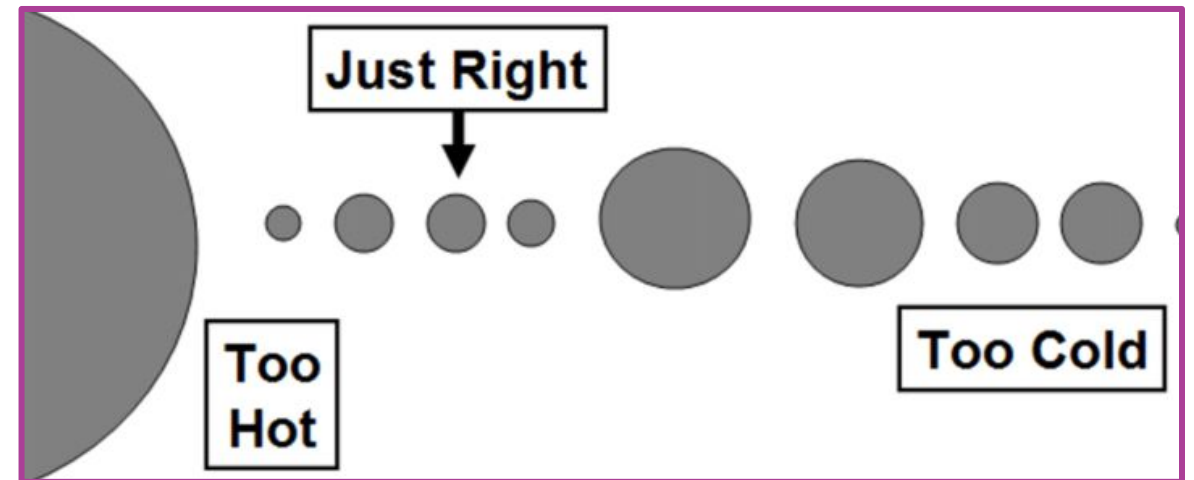
- ✓ terms & classes (cell-based assays)
- ✓ fit for commercial purposes

## 2. Complex

- ✓ Large number of assay terms to cover breadth of offerings from Eurofins Discovery

## 3. Maintenance

- ✓ Currently active and maintained
- ✓ Identified owners

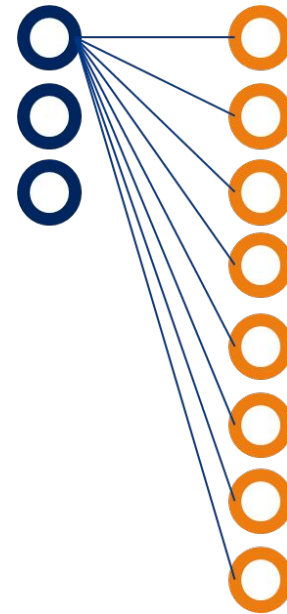


Baum, Seth. (2013). *Journal of Sustainability Education*.

Interesting, but minimally helpful in getting our work done.

Sample terms to match:

- *affinity tag purification*
- *affinity binding assays*
- *inflammation assays*
- *assays for inflammation*



- ~500 terms from one eMarketplace against ~5K terms from another
- ~2.5M pair-wise comparisons by eyeballing
- one pair/second...that's about 30 days non-stop

Needs a more analytical approach.

## Approaches in increasing usefulness

- Eyeball
- Excel VLOOKUP
  - Exact match
  - Approximate match
- NLP
  - Levenshtein distance
  - Cosine similarity









# A Better Approach than Eyeballs?

## Excel VLOOKUP options

- Exact match - matches exact terms well
- Approximate match - only looks at matching characters at the beginning of the words
  - catapult and caterpillar will be flagged as matches!

In cases where there are few exact matches, Excel doesn't help much. What about approximate matches?

 bat bag bar	 rag rat ran
 fan fen fat	 mat mad map
 can cat cap	 had has hat

<https://www.myteachingstation.com/>

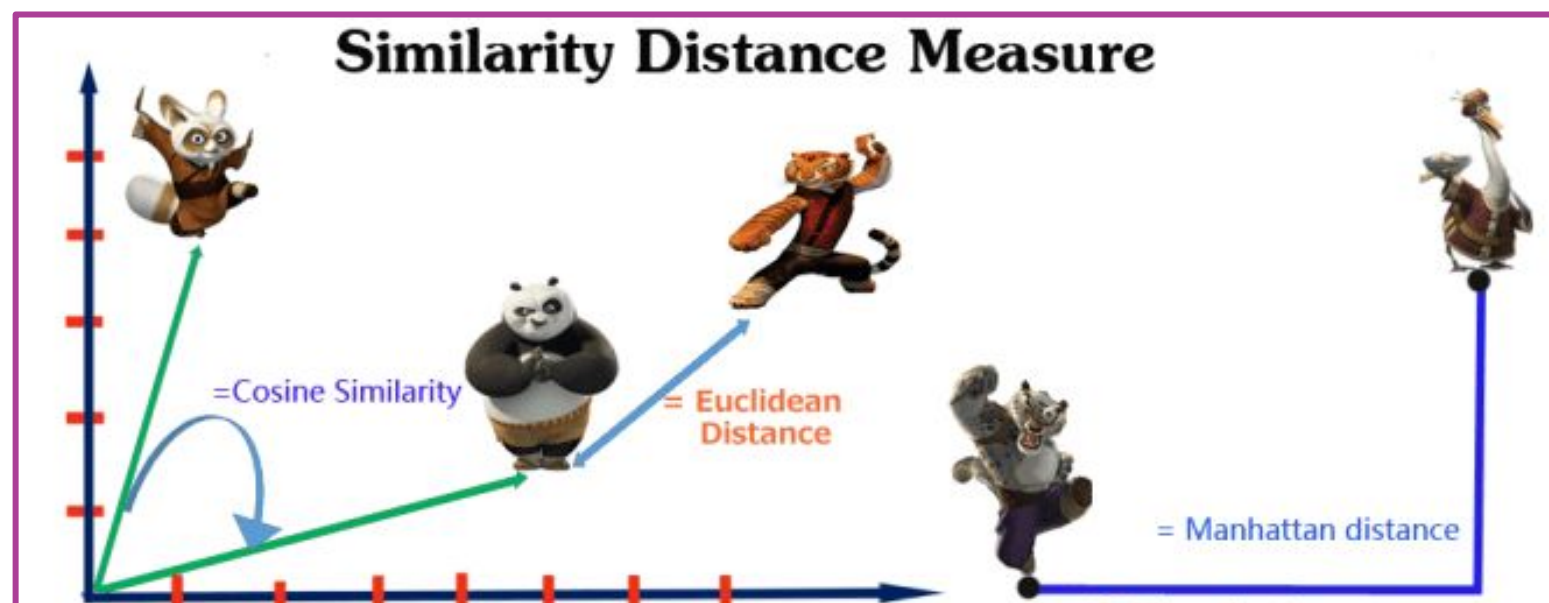
# Approximate Matches Produced by Excel VLOOKUP

Search Term	Approximate Match
Basic Metabolic Profiling	Basic Metabolic Panel
Behavioral Testing Models	Behavioral Phenotyping
Bioequivalence Studies	Biodistribution Studies
BioMAP® Human Phenotypic Platform	Biology
Biomarker Services	Biomarker Discovery
Biomarker Services Consulting	Biomarker Discovery
Biosafety Services	Bioproduct Safety Testing
Biotin Protein Conjugation	Biostatistics & Bioinformatics
Biotin-Antibody Conjugation	Biostatistics & Bioinformatics
Bone Models	Bone Mineral Density Testing

# A Better Approach: Computational String Comparison

- Comparing words, phrases, terms – collectively called strings – is a non-trivial task
- Finding exact matches is easy, but finding a degree of similarity between strings is difficult
- A change in a single character can signal a small or a large difference in meaning:
  - Assay : assays – similar meaning
  - 3d cell culture : 2d cell culture – different meaning
  - Apparently different strings can carry the same meaning:  
human – homo sapiens

- Various metrics have been developed to measure similarity between terms
- We have used two similarity metrics and compared their usefulness
  - Levenshtein distance
  - Cosine similarity



<https://dataaspirant.com/five-most-popular-similarity-measures-implementation-in-python/>

- The minimum number of single-character edits (insertions, deletions or substitutions) required to change one word into the other
- For example, the Levenshtein distance between "kitten" and "sitting" is 3:

1. kitten → sitten (substitute "s" for "k")
2. sitten → sittin (substitute "i" for "e")
3. sittin → sitting (insert "g")



[https://en.wikipedia.org/wiki/Levenshtein\\_distance](https://en.wikipedia.org/wiki/Levenshtein_distance)



# Small Distances Don't Always Mean Closely Related

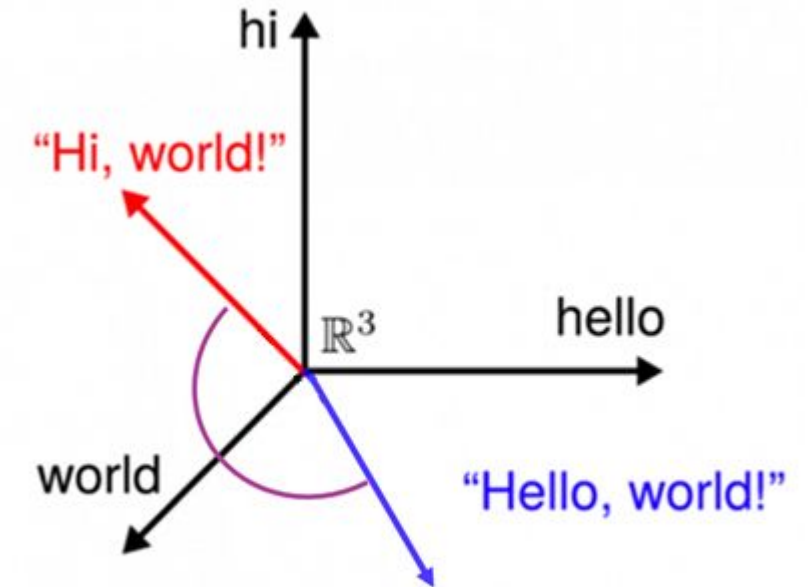
Query Term	Match_1	LSim_1
immunophenotyping	immunophenotyping	0
3d cell culture	2d cell culture	1
caco2 cell line	caco2 cell lines	1
caspase assay	caspase assays	1
cell based assays	cellbased assays	1
gene set enrichment analyses	gene set enrichment analysis	1
human buffy coat	human buffy coats	1
protease assay	protease assays	1
pulmonary hypertension animal model	pulmonary hypertension animal models	1
rna extraction	dna extraction	1
western blots	western blot	1
decarboxylase assays	carboxylase assays	2
facs	aas	2
vitro tme models	vitro eye models	2
sar assays	nab assays	2
aldolase assays	hydrolase assays	3
asthma	asapms	3
bone models	mouse models	3



Levenshtein distance can be misleading

## Second Metric: Cosine Similarity

- Cosine similarity is a measure of similarity between two non-zero vectors of an inner product space
- It is defined to be equal to the cosine of the angle between the vectors
- “S” is measured in the range 0 to 1
  - Exact match if  $S = 1$
  - Approximate match if  $1 > S \geq 0.5$
  - No match if  $S < 0.5$



<https://medium.com/@adriensieg/text-similarities-da019229c894>

- The Python Sklearn package contains algorithms for the necessary code-fu
  - Vectorising
  - Natural Language Processing
  - Similarity calculation
  - Ranking
- String similarity was calculated in Python based on article by Dario Radečić

<https://towardsdatascience.com/calculating-string-similarity-in-python-276e18a7d33a>



[https://en.wikipedia.org/wiki/Cosine\\_similarity](https://en.wikipedia.org/wiki/Cosine_similarity)



# Matches Ranked by Term Similarity: Excel Output

Query Term	Match_1	CSim_1	Match_2	CSim_2	Match_3	CSim_3	Match_4	CSim_4	Match_5	CSim_5	Match_6	CSim_6
Systemic lupus erythematosus (SLE) Human	Systemic Lupus Erythematosus (SLE) Human	0.816	Systemic Lupus Erythematosus (SLE) Human	0.816	Systemic Lupus Erythematosus (SLE) Human	0.816	Systemic Lupus Erythematosus (SLE) Human	0.756	Systemic Lupus Erythematosus (SLE) Human	0.632	Lupus Animal Models	0.289
Cancer Systems Biology	Systems Biology	0.816	Translational Systems Biology	0.816	Translational Systems Biology	0.816	Translational Systems Biology	0.816	Translational Systems Biology	0.816	Translational Systems Biology	0.408
Thermal and Acoustic Imaging	Thermal Imaging	0.816	Thermal Imaging	0.816	Thermal Imaging	0.816	Thermal Imaging	0.816	Thermal Imaging	0.816	Thermal Imaging	0.408
Thrombosis Models	Thrombosis Animal Models	0.816	Thrombosis Models	0.816	Thrombosis Models	0.816	Thrombosis Models	0.816	Thrombosis Models	0.816	Thrombosis Models	0.5
Wistar Rat	Wistar Rat Model	0.816	Wistar Rat Model	0.816	Wistar Rat Model	0.816	Wistar Rat Model	0.816	Wistar Rat Model	0.816	Wistar Rat Model	0.408
Zucker Rat	Zucker Rat Model	0.816	Zucker Rat Model	0.816	Zucker Rat Model	0.816	Zucker Rat Model	0.816	Zucker Rat Model	0.816	Zucker Rat Model	0.408
Inflammatory Pain Models	Acute Inflammatory Pain Animal Models	0.775	Inflammatory Pain Animal Models	0.775	Custom Pain Animal Models	0.577	Neuropathic Pain Animal Models	0.577	Nociceptive Pain Animal Models	0.577	Postoperative Pain Animal Models	0.577
In Vivo Whole Tissue and Animal Imaging	Animal in vivo Imaging	0.775	Whole Brain Imaging	0.671	Animal Whole Tissue Imaging	0.671	Small Animal Imaging	0.6	Animal Imaging	0.516	Animal Tissue Imaging	0.516
Non GLP Large Molecule Bioanalysis	Large Molecule Bioanalysis	0.775	Molecule Bioanalysis	0.516	Bioanalysis	0.447	Bead-Based Bioanalysis	0.316	GLP Audit Bioanalysis	0.316	Plasma Bioanalysis	0.316
Non GLP Small Molecule Bioanalysis	Small Molecule Bioanalysis	0.775	Molecule Bioanalysis	0.516	Bioanalysis	0.447	Bead-Based Bioanalysis	0.316	GLP Audit Bioanalysis	0.316	Plasma Bioanalysis	0.316
Amino Acid Analysis (AAA)	Amino Acid Substitution Analysis	0.75	Amino Acid Analysis	0.75	Hydrolyzed Amino Acid Analysis	0.75	Bile Acid Analysis	0.577	Fatty Acid Analysis	0.577	Nucleic Acid Analysis	0.577
Copy Number Variant Analysis	Gene Copy Number Analysis	0.75	Copy Number Variant Analysis	0.577	Genetic Variant Analysis	0.577	5-Batch Array Analysis	0.354	Analysis of Copy Number Variants	0.354	Biomarker Analysis	0.354
Gene Set Enrichment Analyses	Gene Set Enrichment Analysis	0.75	Gene Set Enrichment Analysis	0.354	Gene Characterization	0.354	Gene Fragmentation	0.354	Gene Knockout	0.354	Gene Synthesis	0.354
Antibody Efficacy Testing In Vivo	In vivo Drug Efficacy Testing	0.75	Drug Efficacy Testing	0.707	Ex vivo Drug Efficacy Testing	0.671	Antimicrobial Efficacy Testing	0.577	In vivo Antibody Efficacy Testing	0.577	In vivo Biosensor Efficacy Testing	0.577
LPS Lung Inflammation Models	Lung Inflammation Animal Models	0.75	Inflammation Models	0.577	Inflammation Models	0.577	Acute Inflammation Models	0.5	Chronic Inflammation Models	0.5	LPS-Induced Inflammation Models	0.5
Rabbit Monoclonal Antibody Development	Monoclonal Antibody (mAb) Development	0.75	Antibody Development	0.707	Therapeutic Antibody Development	0.671	BiTE Antibody Development	0.577	Polyclonal Antibody Development	0.577	Monoclonal Antibody Development	0.5
Pulmonary Hypertension Animal Models	Pulmonary Hypertension Animal Models	0.75	Animal Model	0.577	Diuresis Animal Model	0.577	Hypertension Animal Model	0.577	Animal Model	0.5	Animal Model	0.5
In Vivo Whole Tissue and Animal Imaging	Animal in vivo Imaging	0.707	Whole Brain Imaging	0.612	Animal Whole Tissue Imaging	0.612	Tissue Analysis	0.577	Small Animal Imaging	0.548	AFM Imaging	0.471
Circulating Antibody Assay Development	Antibody Development	0.707	Antibody Development	0.707	Antibody Screening	0.577	Assay Development	0.577	Biochemical Assay Development	0.577	Biomarker Assay Development	0.577
Antibody and Protein Products Handling	Antibody Products	0.707	Protein Sequencing	0.577	Membrane Protein Sequencing	0.5	Alginate Purification	0.354	Antibody Characterization	0.354	Antibody Characterization	0.354
Autoimmune disease	Autoimmune Disease Animal Models	0.707	Models of Disease	0.408	Celiac Disease	0.408	In vitro Disease	0.408	Infectious Disease	0.408	Lyme Disease	0.408
Autoimmune Models	Autoimmune Disease Animal Models	0.707	Central Autoimmune Encephalomyelitis	0.577	Experiment	0.577	Experiment	0.577	Avian Model	0.5	Chicken Model	0.5
BALB/c mouse	BALB/c Inbred Mouse Model	0.707	Mouse Model	0.707	Mouse Model	0.5	Mouse Phenotype	0.5	A/J Mouse	0.408	Agouti Mouse	0.408
Sample Bioanalysis	Bioanalysis	0.707	Based Bioanalysis	0.5	Plasma Bioanalysis	0.5	Sample Preparation	0.5	Sample Preparation	0.5	Clinical Sample Preparation	0.408
Chemokine Biomarkers	Biomarkers	0.707	Biomarkers	0.5	Serum Biomarkers	0.5	Pharmacokinetics	0.408	1,25-Dihydroxy Vitamin D3	0	10X Genomics	0
Cytokine Biomarkers	Biomarkers	0.707	Biomarkers	0.5	Cytokine Analysis	0.5	Serum Biomarkers	0.5	Cytokine Receptor	0.408	Cytokine Signaling	0.408
Translational Biomarkers	Biomarkers	0.707	Biomarkers	0.5	Serum Biomarkers	0.5	Pharmacokinetics	0.408	Translation	0.354	1,25-Dihydroxy Vitamin D3	0
Metastasis Models	Bone Metastasis Animal Models	0.707	Models	0.5	Chicken Model	0.5	Insect Model	0.5	Invertebrate Model	0.5	Mammalian Model	0.5
CRISPR	CRISPR Bioinformatics	0.707	Cloning	0.707	CRISPR Screening	0.707	CRISPR Deletion	0.577	CRISPR Signaling	0.577	Custom CRISPR	0.577
Custom Antibodies	Custom Affinity Purification of Antibodies	0.707	Manufacturing	0.5	Custom Microarray	0.5	Custom Protein	0.5	Llama Antibody	0.5	Custom Antigen	0.408
Model Development	Custom Cell Model Development	0.707	Custom Animal Model Development	0.707	Algorithm Development	0.5	Antibody Development	0.5	App Development	0.5	Assay Development	0.5
Pain Models	Custom Pain Animal Models	0.707	Neuropathic Pain Animal Models	0.707	Nociceptive Pain Animal Models	0.707	Postoperative Pain Animal Models	0.707	Visceral Pain Animal Models	0.707	Acute Inflammation	0.632
Cytotoxicity	Cytotoxicity Assays	0.707	Real-Time Cytotoxicity Assay	0.577	Medical Device	0.5	Natural Killer Cell	0.447	1,25-Dihydroxy Vitamin D3	0	10X Genomics	0
Electroporation Transfection	Electroporation	0.707	Transfection	0.707	Chemical Transfection	0.5	In utero Electroporation	0.5	Microinjection	0.5	Neuronal Electroporation	0.5
Endocrine Models	Endocrine Disease Animal Models	0.707	Avian Models	0.5	Chicken Model	0.5	Insect Model	0.5	Invertebrate Model	0.5	Mammalian Model	0.5



# Cosine Similarity Approximate Matches Work Fairly Well

Query Term	Match_1	CSim_1	Match_2	CSim_2
Acetyltransferase Assays	Acetyltransferase Assays	1	Histone Acetyltransferase (HAT) Assays	0.707
Activity Prediction	Activity Prediction	1	Animal Antiplasmin Activity	0.408
Acute Kidney Failure Models	Acute Kidney Failure Models	1	Acute Liver Failure Animal Models	0.671
wound healing	In vitro Wound Healing Assay	0.707	Wound Healing Animal Models	0.707
In Vivo Efficacy	In vivo Drug Efficacy Testing	0.707	Ex vivo Drug Efficacy Testing	0.632
LCMS Analysis	LC-MS	0.707		
Multidrug Resistance	Multidrug Resistance (MDR) Testing	0.707		
DNA Mutagenesis	Mutagenesis	0.707		
Membrane Permeability (PAMPA)	Parallel Artificial Membrane Permeability Assay	0.707		
Respiratory Models	Respiratory Disease Animal Models	0.707	Avian Models	0.5
Cardiovascular Toxicology	Toxicology	0.707	Computational Toxicology	0.5
Cell Transfection	Transfection	0.707	B Cell Count	0.5
RNA Interference	Vector Based RNA Interference	0.707	Human RNA	0.5
Human Primary Cell Disease Model	Human Primary Cell Isolation	0.671	Animal Models of Disease	0.516
Small Cell Lung Cancer Models	Lung Cancer Animal Models	0.671	Lung Small Cell Xenograft	0.671
Adherent Cell Culture	2D Cell Culture	0.667	Actinomyces Cell Culture	0.667
Method Development and Validation	Analytical Method Development	0.667	Analytical Method Validation	0.667
Disseminated Disease Models	Animal Models of Disease	0.667	In vitro Disease Models	0.667
Annexin V Apoptosis Assay	Annexin V Staining Assay	0.667	1,25-Dihydroxyvitamin D Assay	0.408
Basic Metabolic Profiling	Basic Metabolic Panel	0.667	Compound Profiling	0.408
Brain Tumor Mouse	Brain Tumor Models	0.667	Brain Xenograft	0.408
Caco-2 Cell Line	Caco-2 Cell Lines	0.667	Cell Line Authentication	0.667

In the best match column (Match\_1)

Out of 523 terms:

- ~20% exact matches
- ~60% approximate matches
- ~20% unresolved

with this particular contributor group

Finding and using the correct ontology is challenging.  
Our ultimate goal is to contribute assay data to an ontology that we find helpful.

- Many biomedical ontologies relevant for cellular assays exist
- Excel is excellent for matching exact terms but poorly matches anything beyond that
- Cosine similarity for string matching works well for approximate matches
- People are still needed for curation
  - A key term, phenotypic screening, missing from the list of 5K terms was added
  - Irrelevant terms (12 %) included by the similarity scoring method were removed

# Could We Have Used A Different Ontology Mapping Tool?

- Short time frame and resources to complete the project
- We have little expertise with ontologies
- Although there are free ontology mapping tools online...  
...the terms we worked with - we think - were not encompassed by any ontology

# Thank You

Aimee Allen  
Jennifer Drake  
Ellen Berg

- Categories & terms were equally understood / meant the same thing to contributors of diverse backgrounds?
- eMarketplaces were educated about ontologies & established adherence to controlled ontologies?
- Terms were harmonized across content platforms beyond eMarketplaces?
  - assay descriptions
  - web content
  - file properties and keyword tags

