45

# Metadata vs Metadata: Aiming for the Translation of Metadata Across Different Omics

吉沢明康、高橋悠志、石野公基、奥田修二郎

Akiyasu C. Yoshizawa, Yushi Takahashi, Koki Ishino, Shujiro Okuda 新潟大・医・メディカル AI センター

Medical Al Center, Niigata University Medical School

### 背景と目的

データリポジトリの意義は、例えばメタ解析、即ち多数のデータの解析結果から全体の傾向を見出す 解析のような、データを再利用する研究に資することである。公開データから各オミクスのメタ解析を 行い、その結果を統合する"疑似トランスオミクス解析"研究は今後の増加が予想され、それに備える ことはまた、データリポジトリの今後の課題と考えられる。

しかし実際には、メタデータの項目・用いられている統制語彙などはオミクスによって異なり統一 されていない。メタデータの内容を規定する基準は各分野に存在するにもかかわらず、具体的な語彙の 使用については規程が存在しないため、同一概念に対して使用されている語彙がオミクス分野によって 異なる場合さえある。

我々は今までに、プロテオーム・データリポジトリ&データベース jPOST 及びメタボローム・データ 解析パイプライン Shin-MassBank の 2 つのプロジェクトに於いてデータリポジトリを開発し、メタデータ の処理を実装してきた。一般にプロテオームとメタボロームは、データ処理のプロセスが質量分析の応用 分野の中では最も類似しているとされるが、実際に jPOST のメタデータをベースに Shin-MassBank の リポジトリ MB-POST 用のメタデータを作成する際に比較したところ、データの種類・データの構造・ 実際に用いられる語彙など多くの点に差異が見られた。

本発表では、将来的なリポジトリ・データ解析の統合を見据えて、これら両分野のメタデータの対応 づけを開始した。

### プロテオーム対メタボローム:メタデータ比較(1)

#### メタデータの構成

・最少情報基準のあるプロテオミクスでは、**生データファイル毎**にそのデータの属性(試料の 内容や行った前処理など)を記述する**SDRF** (Sample Data Relation Format) 情報が用いられる。

1	source name	characteristics[organism]	characteristics[organism part]	characteristics[cell line]	characteristics[ancestry catego		コーノルタバミュギナムアルフ			
2	Sample 1	Homo sapiens	not applicable	HeLa	Black		ファイル名が記載されている			
3	Sample 10	Homo sapiens	not applicable	proteomic profiling by mass spectrometry  proteomic profiling by mass spectrometry		https://ftp.pride.ebi.ac.	uk/pride/data/archive/2014/07/PXD000396/120315QEx2_RS1_50nl-min_10ngHeLa_1h_01.raw	1		
4	Sample 11	Homo sapiens	not applicable							
5	Sample 12	Homo sapiens	not applicable	proteomic profiling by mass spectrometry		https://ftp.pride.ebi.ac.uk/pride/data/archive/2014/07/PXD000396/120323QEx2_RS1_20nl-min_0k1HeLa_10h_01.raw				
6	Sample 13	Homo sapiens	not applicable	proteomic profiling by mass spectrometry		https://ftp.pride.ebi.ac.uk/pride/data/archive/2014/07/PXD000396/120315QEx2_RS1_50nl-min_10ngHeLa_1h_02.raw				
7	Sample 14	Homo sapiens	not applicable	proteomic profiling by mass spectrometry		https://ftp.pride.ebi.ac.uk/pride/data/archive/2014/07/PXD000396/120323QEx2_RS1_20nl-min_0k1HeLa_5h_01.raw				
8	Sample 15	Homo sapiens	not applicable	proteomic profiling by mass spectrometry		https://ftp.pride.ebi.ac.uk/pride/data/archive/2014/07/PXD000396/120323QEx2_RS1_20nl-min_0k1HeLa_3h_01.raw				
9	Sample 16	Homo sapiens	not applicable	proteomic profiling by mass s			uk/pride/data/archive/2014/07/PXD000396/120315QEx2_RS1_20nl-min_10HeLa_8h_01.raw	1		
10	Sample 17	Homo sapiens	not applicable	proteomic profiling by mass s			uk/pride/data/archive/2014/07/PXD000396/120330QEx2_RS1_50nl-min_100ngHeLa_14h_01.raw	1		
	PXD000396			proteomic profiling by mass s			uk/pride/data/archive/2014/07/PXD000396/120315QEx2_RS1_50nl-min_10ngHeLa_5h_01.raw			
	https	://aithub.com/biabio/i		proteomic profiling by mass s adata/blob/8297770			uk/pride/data/archive/2014/07/PXD000396/120309QEx2_RS1_50nl-min_0k1HeLa_3h_01.raw 4fd327695b8/annotated-projects/PXD000396/PXD000396.sdrf.tsv	1		

・これに対してメタボロミクスでは、ISA (Investigation, Study and Assay) データモデルに基づい たメタデータが作成されており、SDRFに類似する情報としては"Study Design"情報が用いられる が、ファイル名は必須ではなく、記載されていないことが多い。

Source Name	Characteristics[Organism]	Term Source REF	Term Accession Number	Characteristics[Organism part]	Term Source REF	Term Accession Nu
Wild Type	Escherichia coli	NCBITaxon	http://purl.obolibrary.org/obo/NCBITaxon_562	Cell	SLSO	http://www.w3.org/20
V610F	Escherichia coli	NCBITaxon	http://purl.obolibrary.org/obo/NCBITaxon_562	Cell	SI SO	httn://www.w3.ora/20
Wild Type	Escherichia coli	NCBITaxon	http://purl.obolibrary.org/obo/NCBITaxon_562	membrane	All Proje	ect Subject Stu
V610F	Escherichia coli	NCBITaxon	http://purl.obolibrary.org/obo/NCBITaxon_562	membrane		

https://www.ebi.ac.uk/metabolights/ws/studies/MTBLS13050/download?file=s MTBLS13050.txt

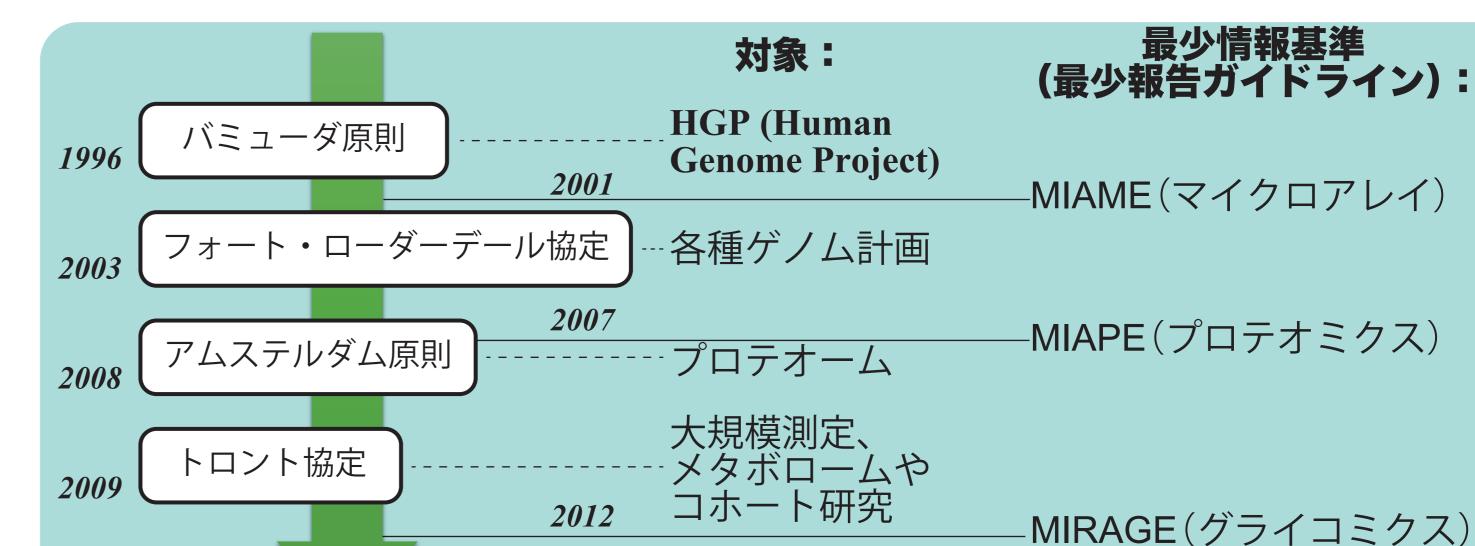
MetaboLights: MTBLS13050

- Study Designで示されるのは基本的に「試料についての 属性」であるが、一つの試料が必ず一つのファイルに記録 されるわけではないので、このデータを再利用(再解析) するためには、他のメタデータと突き合わせることが 最低限必要になる。
- ・この意味で、データ再利用には追加の処理が必要になる。

Metabolomics Workbench: ST004161

# 背景:生データの公開

・測定データ(生データ)を公開する、という制度はヒトゲノム計画のデータを対象に 開始され、順次他のオミクスに拡大した【1】。



- ・この結果、「データは解析されなくても、存在している(産生された)だけで価値が ある」という思想が一般化し、データの再利用(プロテオームなど FDR 制御に Target-Decoy 法を使っている場合はデータの再解析)が普及した。
- ・再利用のためにはデータの属性を説明するメタデータ (metadata) が必須になり、関係 学会によって順次、「最低でもこれだけの情報は必要」という「**最少情報基準** (minimum information standard)」が決定された。但し、例えばメタボローム分野では Metabolome Society が対応できなかったため、メタボロミクスの最少情報基準は存在していない。
- ・このため、特にメタボロミクスでは「基準に寄せる」形で用語を標準化することが難 しい。

### プロテオーム対メタボローム:メタデータ比較(2)

#### メタデータの項目

• Shin-MassBankのリポジトリMB-POST 構築時には、先行するMetaboLights/ Metabolomics Workbench/MetaboBank のメタデータを調査し、それをjPOST (正確にはデータジャーナルJPDM[2]) のメタデータにマップし、更にメタボロ ーム研究者が吟味することによって項目を 決定した。

右パネルに示すように、プロテオミクス では測定方法・ソフトウェア解析の段取り が固まっているのに対し、メタボロミクス では質量分析の測定方法がより複雑である。

- 例えばプロテオミクスでは陽イオンモードでの 測定のみであるのが通常であるが、メタボロ ミクスでは陰イオンモードの測定も行うことが
- またメタボロミクスでは情報解析の定番 方法論がないため、自由記述の比率が高い。
- ・これらの差異については各分野の測定・ 解析に基づく必然性があるため、無理な

プロテオームとメタボロームのメタデータで 対応するカテゴリにおいてそのオーム特有の内容

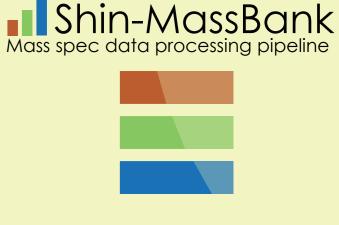


Fractionation 細胞・タンパク質 ペプチドの分離

Enzyme\_mod (消化酵素・ 翻訳後修飾の設定)

MS mode Quantification platform

Software setting (より詳細)



Preparation 誘導体化 内部標準の有無・種類 (対応カテゴリなし)

Analytical\_condition ionization polarity collision energy

Software\_setting (自由記述多し)

用語の統一よりも、項目を増加させてカバーする(例えばプロテオミクスでは「ionization」に 「ESI」、「polarity」に「positive」、「collision energy」欄は空白、など)のが適切と思われる。

# プロテオーム対メタボローム:メタデータ比較(3)

#### メタデータの語彙

- ・項目同士に対応がついた場合でも、用いられている 語彙自体が異なることがあり、その対応づけが必要に なる。現状では、LLM (Large Language Model) を用い るのが最も適当と考えられる。
- ・ 予備調査として、「生物種名と学名の対応」について、 各種の LLM が持つ情報について調査した。
- LLM: multilingual-e5-large (e5), all-MiniLM-L6-v2 (l6), Meta-Llama-3.1-8B-Instruct (llama)

各 term を LLM で embedding vector に変換、各ベクトル間の (1-cosine 類似度)を計算、特徴空間内での term 間距離を得る。 なお llama はプロンプトからの質問「human の学名は何ですか?」 「学名が Homo sapiens の生物種は何ですか?」に正しく返答している

・右表に示すように、単語間の"距離"を単純に比較す る場合に、一般名と学名をタイトに結びつけている LLN は、調査したものにはなかった。但し I6 は「ヒト→学名 と「ウェスタン・ローランド・ゴリラとその学名の双 方向」を正しく第1位にしており、一般名と学名の対応 づけを目指す追加学習には適切かもしれない。

### 各種言語モデル中の生物種名・学名の類似度の例

Instrument QC

Instrument QC

QC(I)#7

QC(I)#5

	label	n1	n2	n3	
	Homo sapiens	Caenorhabditis elegans	chimpanzee	Aedes aegypti	
	human	dog	horse	apple	
	Pan troglodytes	Pan paniscus	Bufo gargarizans	Theropithecus gelada	
chimpanzee		cheetah	Plasmodium malariae	bonobo	
	Pan paniscus	Pan troglodytes	Panicum virgatum	Panicum hallii	
bonobo		dingo	chimpanzee	cheetah	
Gorilla gorilla		Marmota marmota marmota	Bison bison bison	Conger conger	
	western lowland gorilla	western yellowjacket	western mosquitofish	western predatory mite	
-	16				
Н	label	n1	n2	n3	
	Homo sapiens	chimpanzee	human	Bornean orangutan	
	human	Homo sapiens	dog	chimpanzee	
	Pan troglodytes	Pan paniscus	Atta cephalotes	Trachymyrmex cornetzi	
	chimpanzee	Gorilla gorilla	pig-tailed macaque	golden snub-nosed monkey	
	Pan paniscus	Pan troglodytes	Panulirus ornatus	Suncus etruscus	
	bonobo	Bolivian squirrel monkey	Gorilla gorilla	rhesus monkey	
	Gorilla gorilla	western lowland gorilla	chimpanzee	green monkey	
	western lowland gorilla	Gorilla gorilla	Plasmodium sp. gorilla clade G2	African savanna elephant	
	llama label	-1	-2	-2	
			n2	n3	
ı	Homo sapiens	Hippocampus comes	Hydra vulgaris	Hyaena hyaena	
L	human		rat	pig	
	Pan troglodytes	Pan paniscus	Theropithecus gelada	Panthera pardus	
	chimpanzee	chum salmon	chicken	cherry salmon	
	Pan paniscus	Panthera uncia	Panthera pardus	Panthera leo	
	bonobo	dingo	wapiti	bilby	
	Gorilla gorilla	Gekko japonicus	Hydra vulgaris	Aquarana catesbeiana	
	and the second and all and a second	l a a a t a sur a la sur a sur a la a	and the same of th	and the same and a selection	

### 【現時点でのまとめ

- ・データの統一的な取り扱いやトランスオミクス解析のために、メタデータ の対応づけは必須である。そのためには最少情報基準の突き合わせが効率 的であるが、基準のない分野のためには"自前"で対応づけを検討する 必要がある。
- ・基準のないメタボロミクスなどの場合、リポジトリ毎にメタデータを調査 し、どのファイルにどの項目が含まれるか調査する必要がある。更に分野 特有のメタデータについては、それらを網羅してメタデータ項目の合併 集合を作成する必要がある。
- ・最も困難があるのは語彙自体の対応づけである。これは「全場合に対応 できる網羅的なオントロジー」の方向よりも、文脈に応じて詳細に意味を 指定する方が現実的と考えられ、そのためには LLM を用いるのが有効と 考えられる。但し LLM はモデルによって内部に持つ情報の構造化に差が あると考えられるため、最適なモデルを選び、その上で目的に特化した 追加学習を行うことが必要と考えられる。
- ・今回の簡単な予備調査でも、プロンプトからの質問への回答が、モデル 内部の構造化情報をそのまま出力しているのではないことが強く示唆され た。適切なモデルの選択・適切な追加学習・適切なプロンプトエンジニア リングの組み合わせが、メタデータの対応づけという問題解決に必要であ ると考えられる。

文献

[1] 吉沢明康・小林大樹・河野信, Proteome Letters, 10, 11-26 (2025)

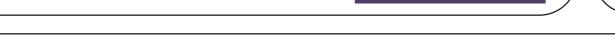
©2025 吉沢 明康 (新潟大学)

[2] 吉沢明康 他, トーゴーの日シンポジウム2024、ポスター16 (2024)



jPOST 及び Shin-MassBank: MB-POST の研究・開発は、科学技術振興機構 (JST)・NBDC 事業推進室による統合化推進プログラム予算に よって実施した(課題番号:JPMJND2304, JPMJND2305)。JPDM の発行には、日本学術振興会の科学研究費補助金・研究成果公開促進 費(国際情報発信強化(B))の支援を受けた(採択課題番号:21HP2004)。

また本研究の実施には、日本学術振興会の科学研究費基金・基盤研究 (B) の支援を受けた (課題番号: 25K03216)。



トーゴーの日シンポジウム 2025

科研費