

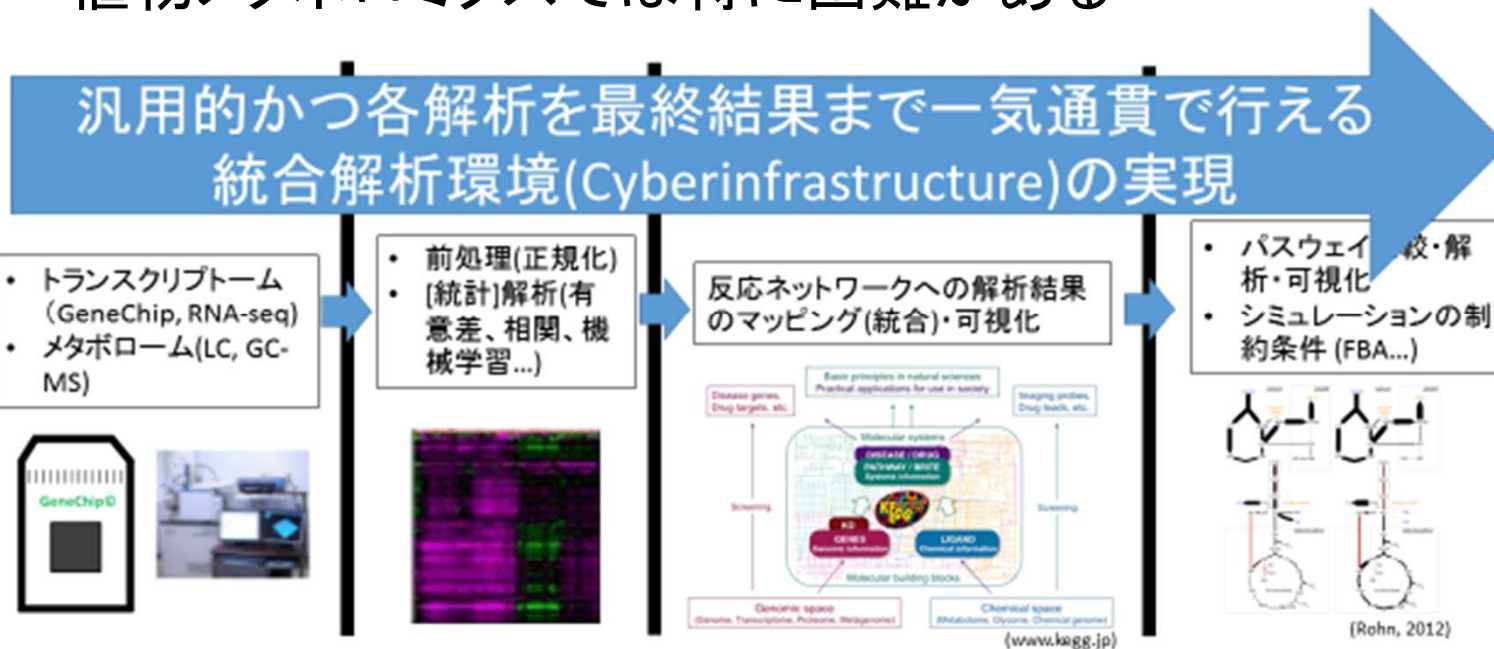
生化学反応ネットワーク統合 解析環境の拡張

理化学研究所
生命システム研究センター
西田孝三

1. 研究概要
2. ユースケースと成果
3. 新展開と今後

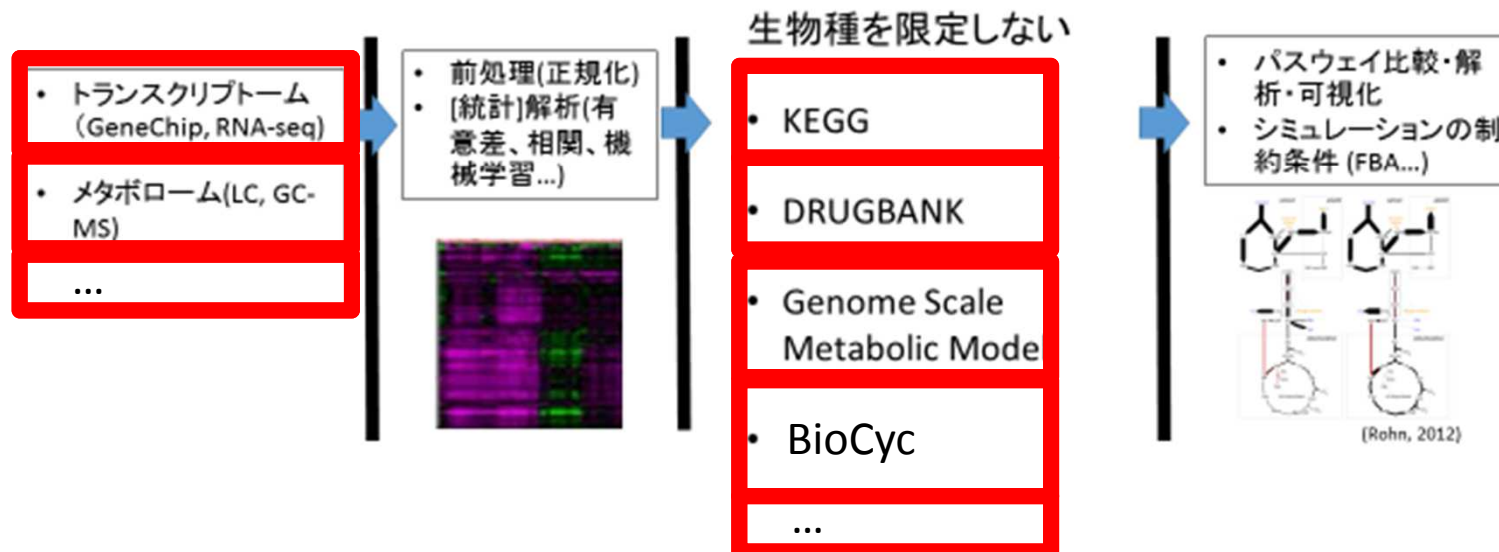
背景、目的

- オミクス実験結果からパスウェイ解析まで一通り行える**フリー**な統合解析環境 (Cyberinfrastructure)の実現
 - 統合化推進プログラムのデータベースを対象
 - 各異種解析を統合する環境
 - 植物メタボロミクスでは特に困難がある



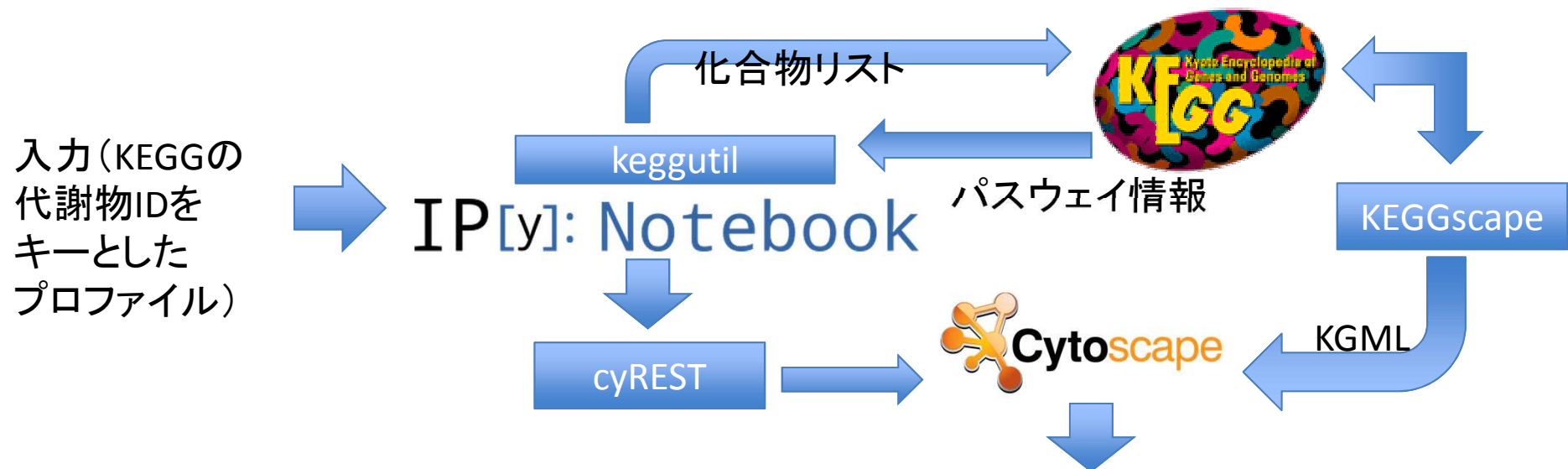
前年度との違い

- 入力情報の拡張（発現プロファイルから代謝物プロファイル）
 - Pathview(Luo, 2013) でも未熟
- 再現性が高く、拡張の用意な解析パイプラインの提供
 - 解析の変更、追加も容易に(リバイズ時など)



研究開発成果概要

- 代謝物プロフィール可視化パイプライン
 - KEGGマッピング用PythonパッケージkeggutilとCytoscape app (cyREST, KEGGscape)を組み合わせたIPython Notebookパイプライン (github.com/kozy2/keggutil, [togotrial2014](https://github.com/togotrial2014))
 - なぜGalaxy(Goecks, 2010)を使わなかったか



出力 (再現性の高い [既存の解析を完全な形で残した] 可視化)

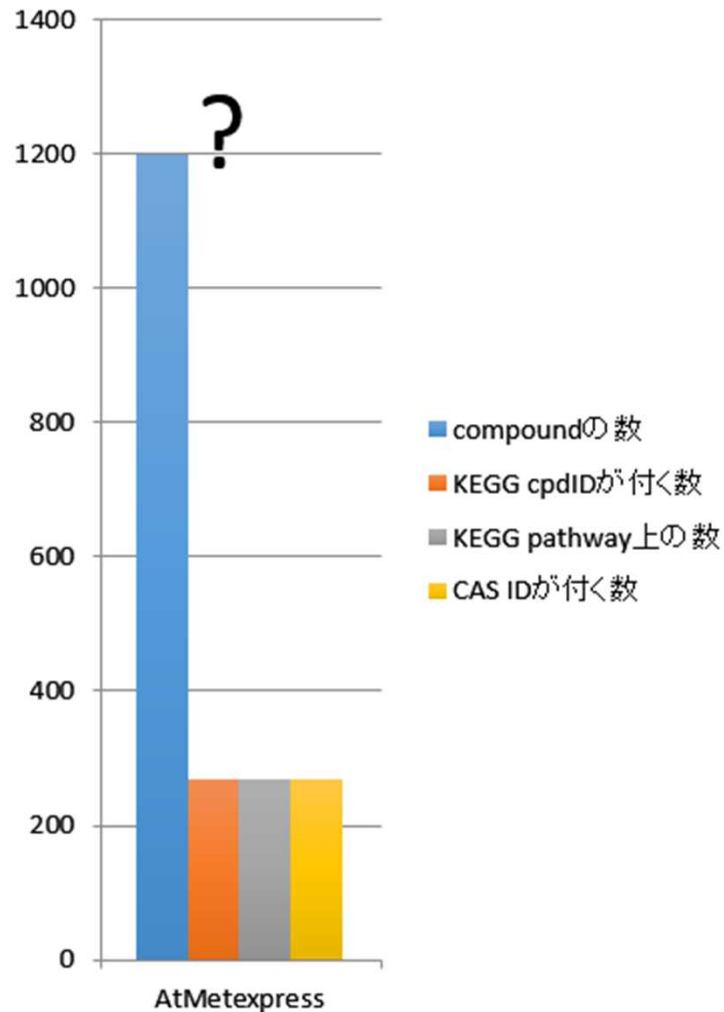
ユースケース

- AtMetExpress
 - 前年度統合解析トライアルの成果DB
 - シロイヌナズナの代謝物~1,200個のプロファイルデータベース
 - パスウェイ統合は未達成
 - データセット名 Matsuda, 2010のみマルチオミクス解析可能な発現プロファイル有
 - 実験条件が合うプロテオミクスデータセット発見できず

結果

- Col0-timecourseのleaf-root間時系列プロファイル比較
 - <http://nbviewer.ipynb.org/github/kozy2/togotrial2014/blob/master/leaf-root.ipynb>
 - <https://github.com/kozy2/togotrial2014/blob/master/leaf-root.svg>

達成できたこと、できなかったこと



できたこと

- 再現性が高く、詳細な化合物プロファイルの可視化

できなかったこと

- 多くの化合物のDB統合
 - 化合物オントロジーを用いたエンリッチメント解析
 - ゲノムスケール代謝モデルとの比較
- パイプラインの汎用化、RDFの活用、データフォーマットの定義

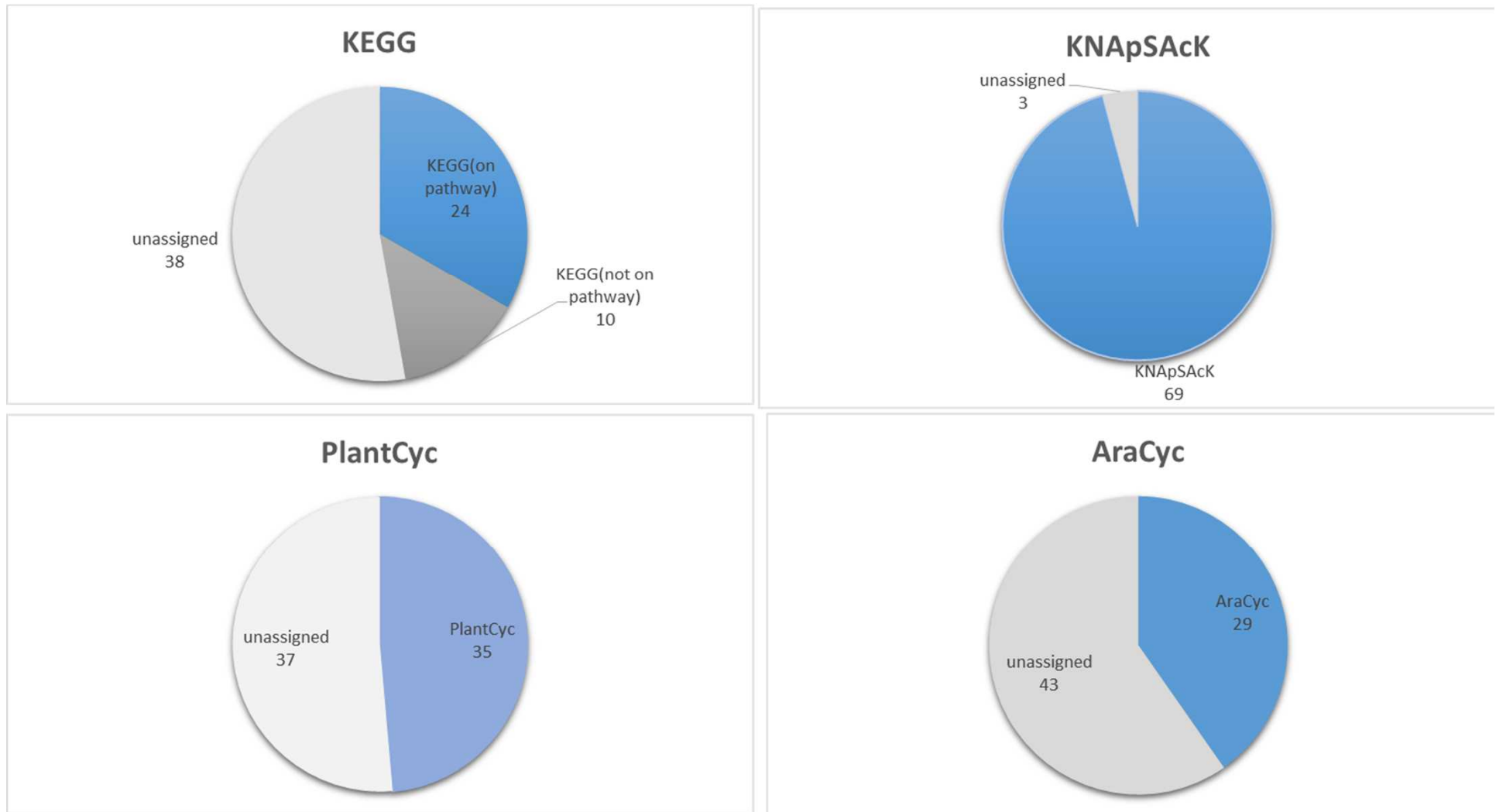
AtMetExpressの パスウェイマッピングにおける難題

- 化合物におけるIDの問題
 - 代謝物のアノテーションのレベルの差
 - 機器の違い、標品の有無、立体化学など
 - 曖昧な化合物名がIDとして用いざるを得なくなる
 - 植物特有の事情(二次代謝パスウェイの重要性)
 - エントリ自体パスウェイデータベースには無いことがある
- マニュアルキュレーションが必要
 - 現状の情報基盤では機械的処理は困難

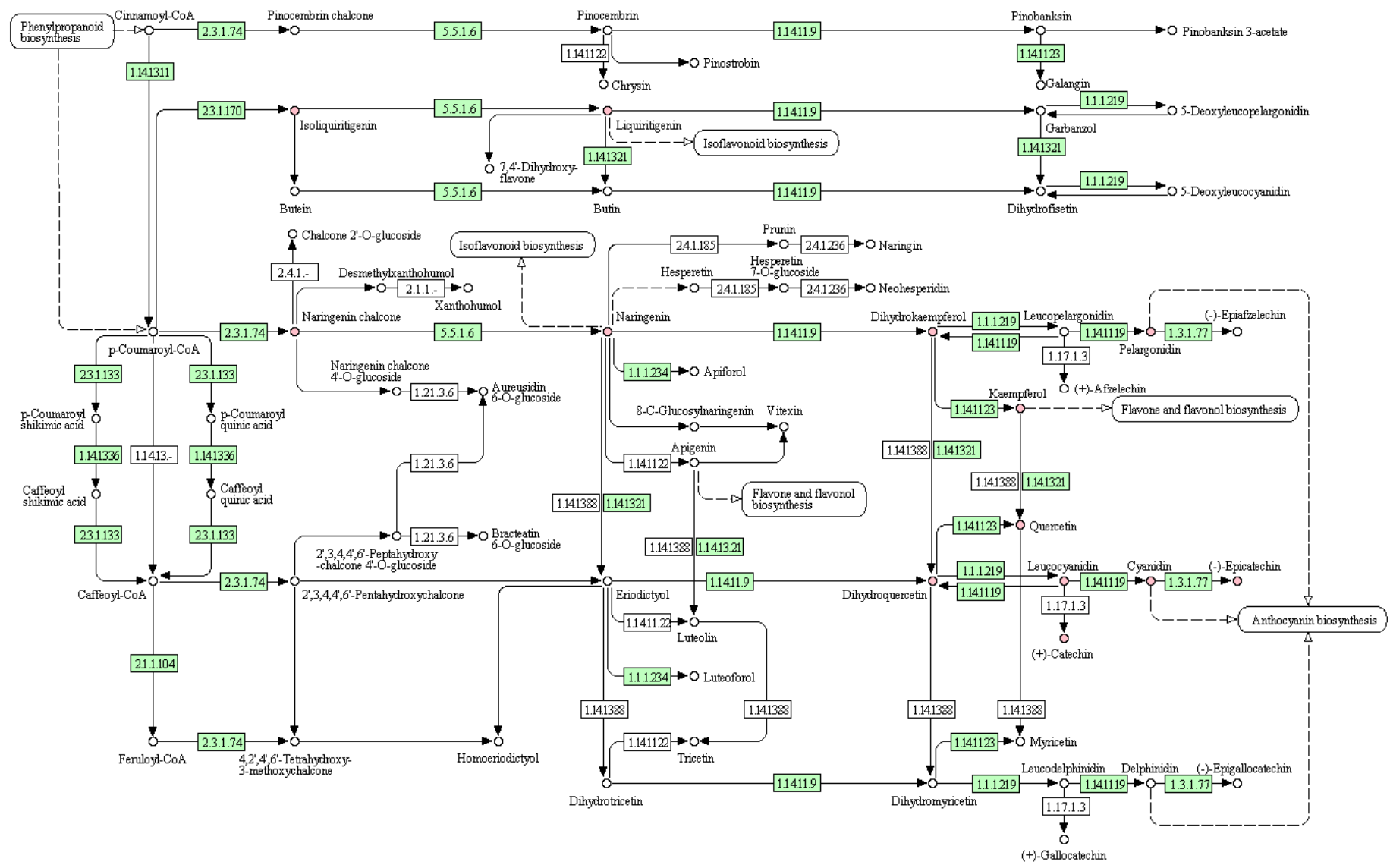
新展開から生まれた開発計画

- 代謝物のマニュアルキュレーション(DBCLS 時松氏に拠る)
 - 二次代謝(Flavonoid)にフォーカス
 - 天然物や二次代謝に特化したKNAPSAcKデータベースを中心に複数DBと照合
 - Saito, Plant Physiol Biochem. 2013 のscaffold structureに基づいたパスウェイ図の作成

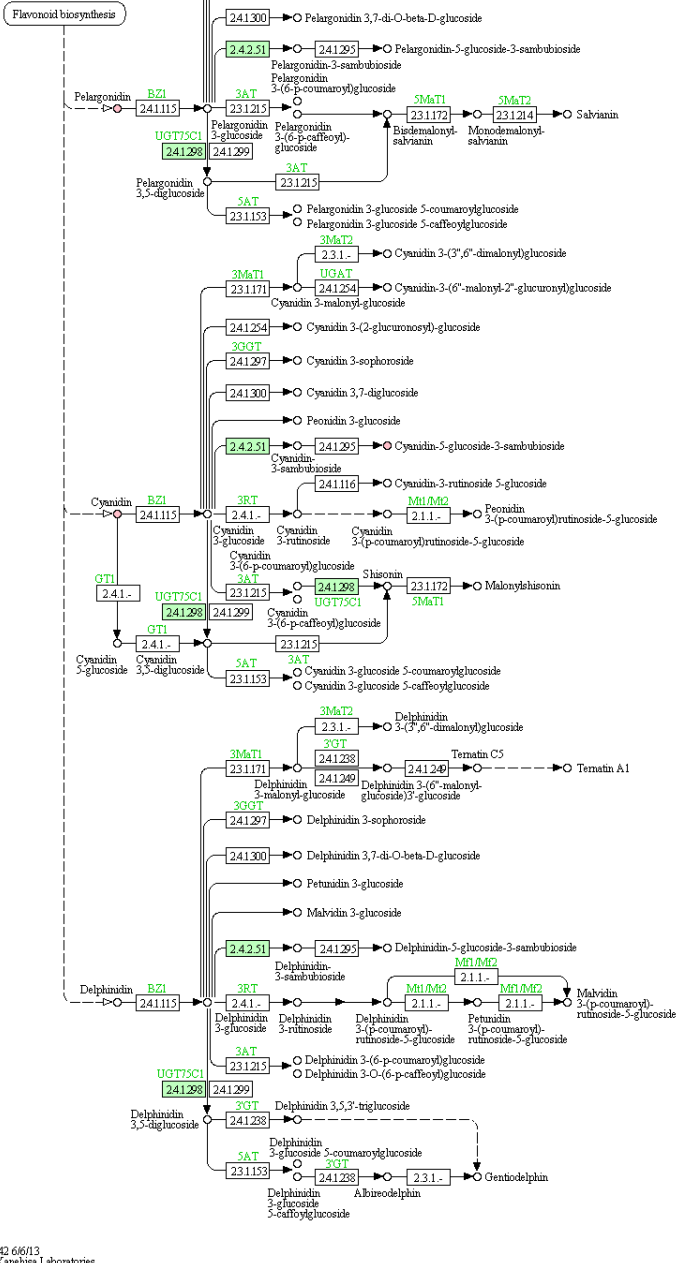
72 Flavonoidsのデータベース 参照可否の図



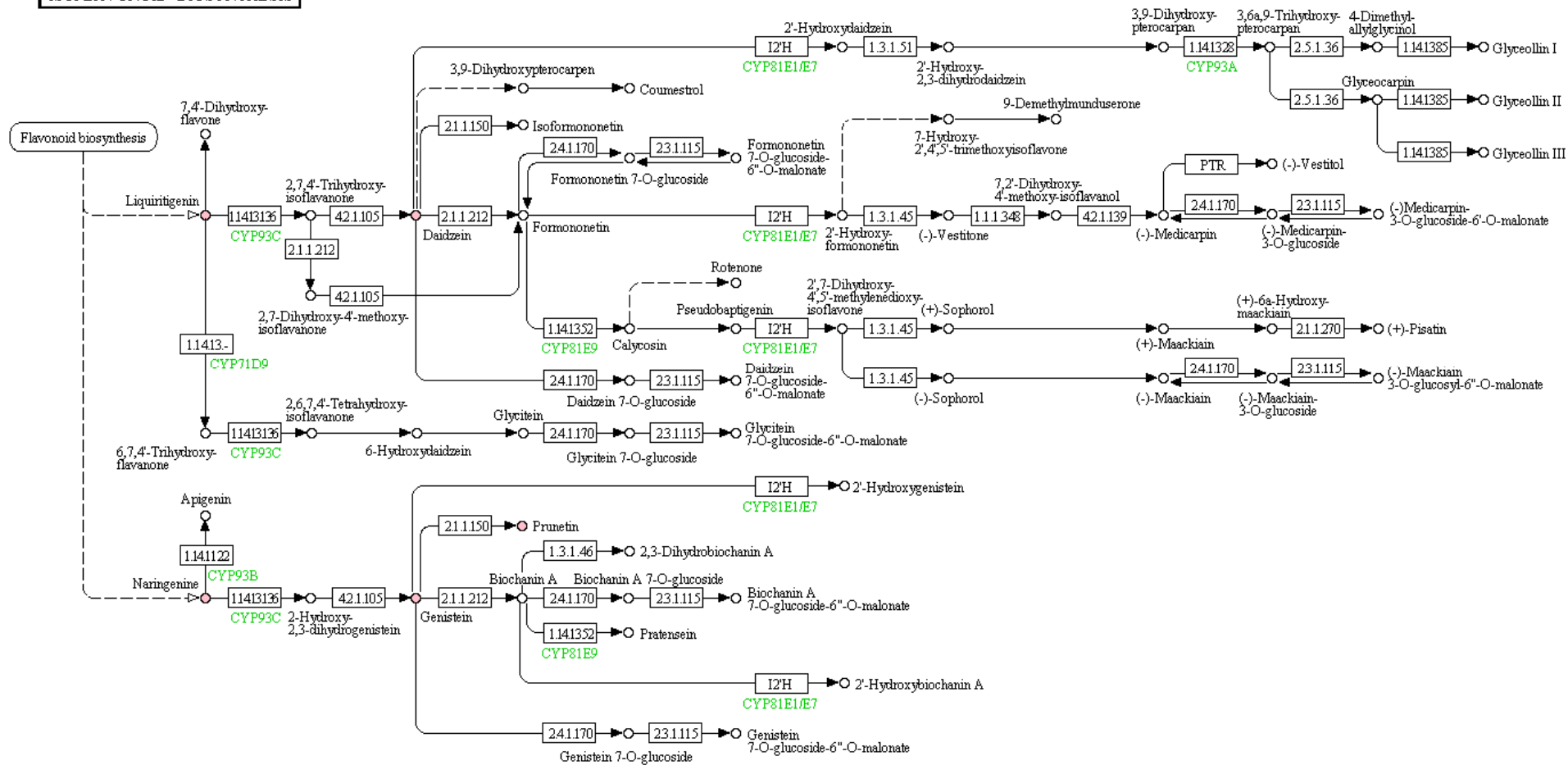
FLAVONOID BIOSYNTHESIS



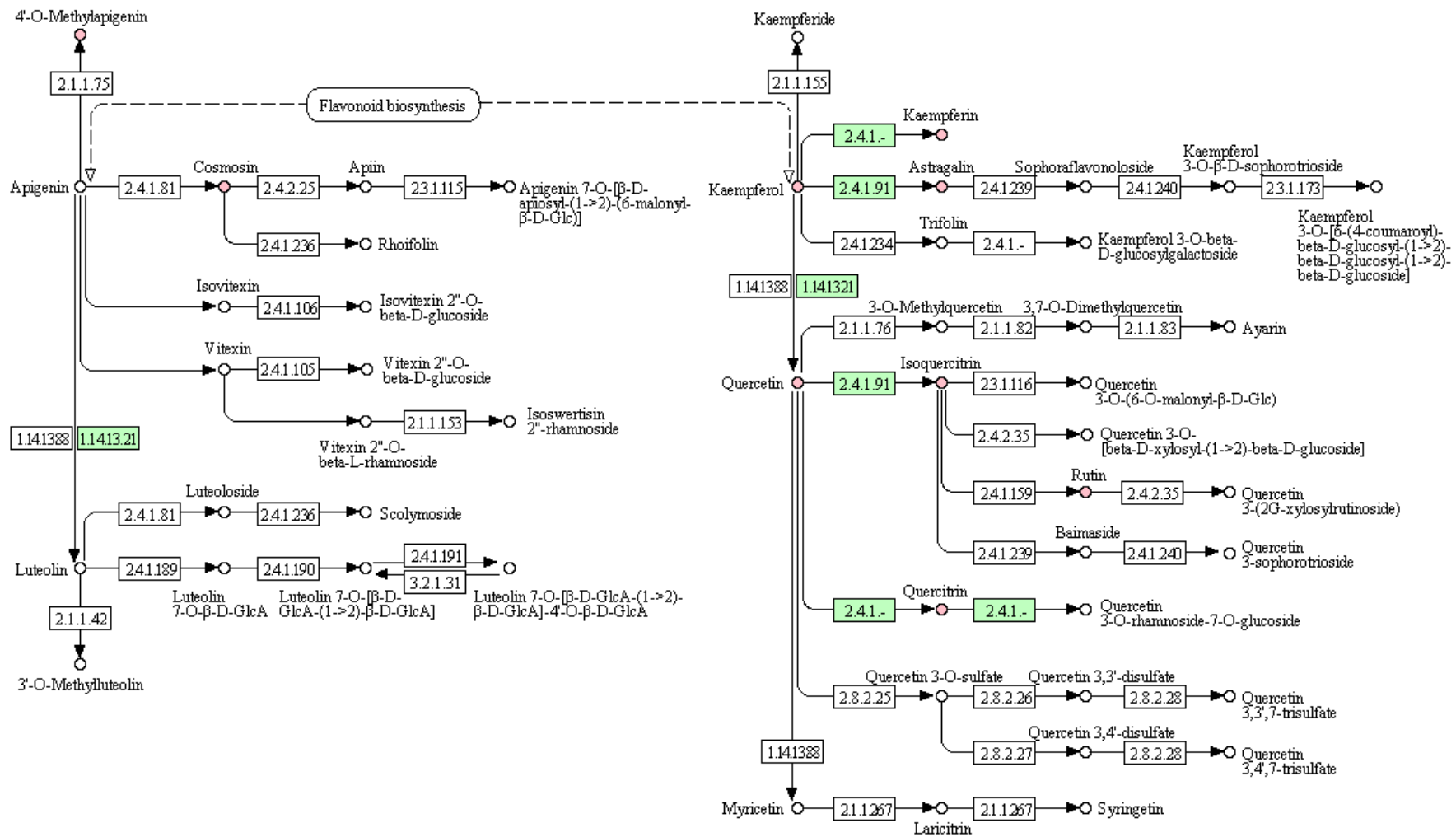
ANTHOCYANIN BIOSYNTHESIS



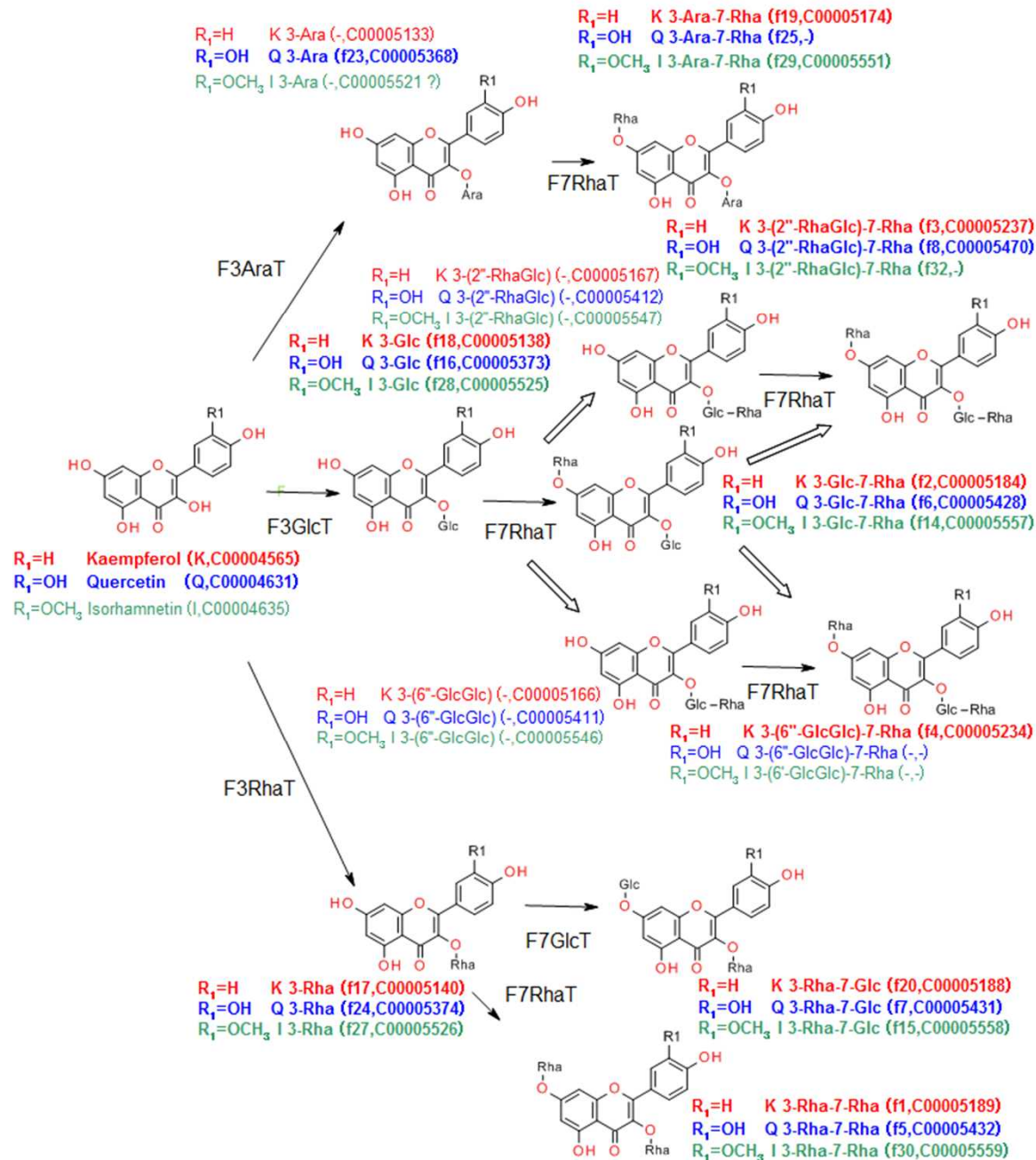
ISOFLAVONOID BIOSYNTHESIS



FLAVONE AND FLAVONOL BIOSYNTHESIS



$R_1=H$ Kaempferol (K)
 $R_1=OH$ Quercetin (Q)
 $R_1=OCH_3$ Isorhamnetin (I)



原子の由来を意識した
 パスウェイ図
 (DBCLS時松氏に拠る)

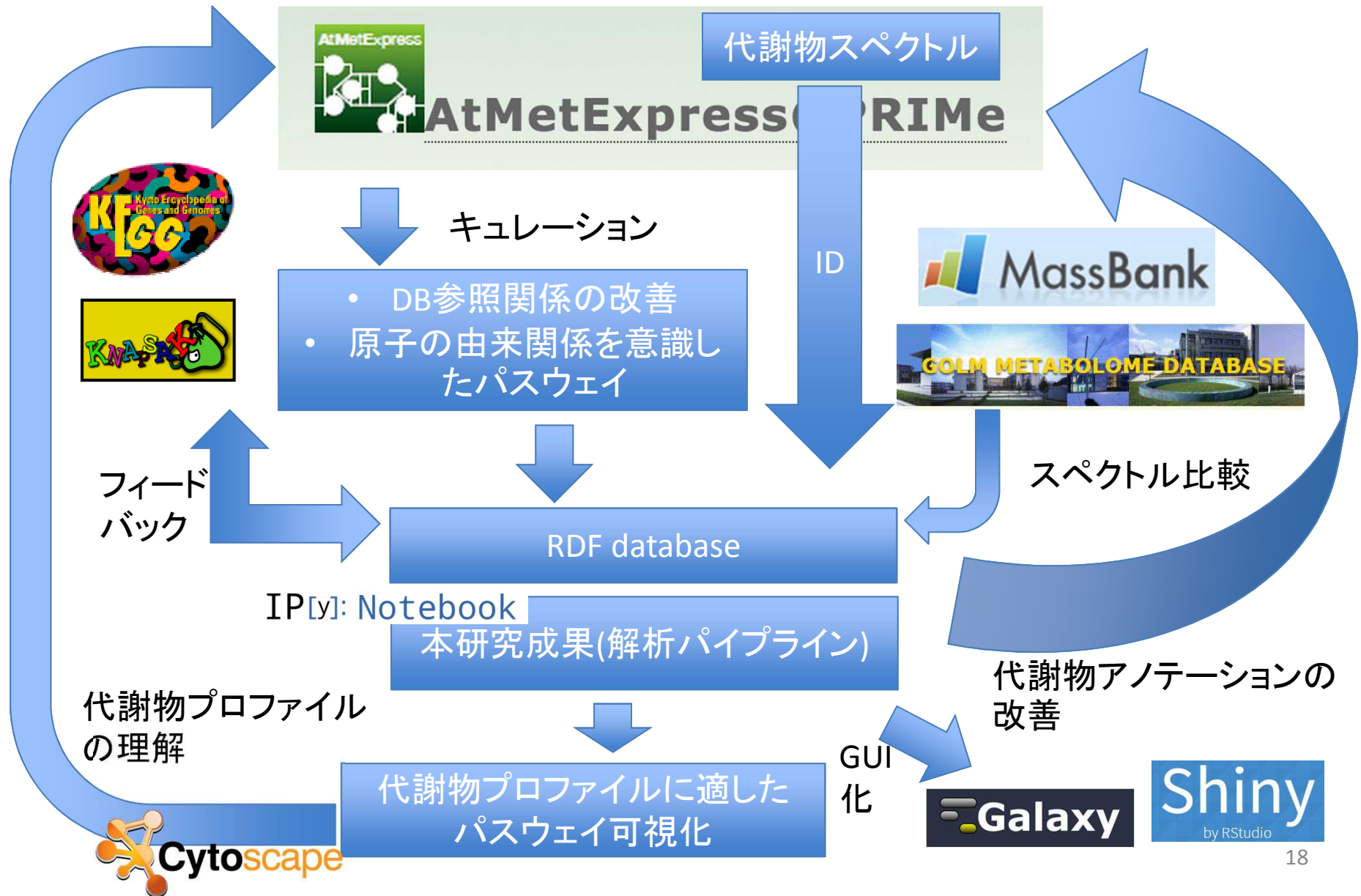
開発の今後

- Glucosinolate, alkaloid, lipidへのキュレーション、パスウェイ図作成(必要に応じ)
- スペクトル情報のリンク、活用
 - Massbank, Golm metabolome database
- RDFデータベース化、公開と解析パイプラインの統合

なぜ代謝物プロフィール の活用にRDFが必要か

- DB参照情報がプロフィールDBで提供できていない
 - 組織を超えた情報リンクを推進する必要がある。
 - データ共有が発現プロフィールより難しい。遅れている
- 自動処理に統一的なデータ保存、取得方法が必要
 - 事前にフォーマットを決める必要がない。代謝物情報を整理するためのフォーマットを定めることは難しい
 - 独自フォーマットはデータ共有を困難にする
- アノテーションの明確化、ID問題解決への手がかかり
 - 曖昧な化合物名への依存からの脱却。スペクトル情報をURIとしアノテーションを定義するなど
 - 単一の情報源には依ることができない。複数の情報源を参照する必要がある

本研究の将来性



まとめ

- 本研究成果パイプラインは代謝物プロファイルの実験間の詳細な可視化を可能にする(時系列など)
 - しかしながらメタボロームにおいてはデータベース統合環境の改善が必要
- 代謝物アノテーション改善のためのRDFを用いたデータベース統合の必要性を示した
 - 原子の由来に基づいた生合成経路オントロジーの土台を提示

謝辞(敬称略)

- 福島敦史 (理研CSRS)
- 時松敏明 (DBCLS)

- 大野圭一郎 (UCSD)
- 瀬々潤 (AIST)

- NBDC/DBCLS biohackathon