# Accurate and Comprehensive Mapping of Multi-omic Data to Biological Pathways

## Application Note

Integrated Biology

## Authors

Anupama Rajan Bhat and Pramila Tata
Strand Life Sciences
Bangalore, India

Stephen Madden
Agilent Technologies, Inc.
Santa Clara, California

## Abstract

This application note describes the use of Agilent-BridgeDB, an essential technology in Agilent's GeneSpring/Mass Profiler Professional (MPP) product to accurately map biological entities on pathways. It describes four case studies that demonstrate how Agilent-BridgeDB enables significantly more accurate mappings between experimentally identified biological entities (for example, genes, metabolites) and the corresponding entities in pathway databases. Common bioinformatics challenges like missing annotations, resolving enantiomers, and incomplete databases are overcome using the Agilent-BridgeDB technology.

**Agilent Technologies**

## Introduction

Pathway analysis provides a useful biological context for differentially expressed entities resulting from the analysis of high-throughput data in any 'omics' (for example genomics, transcriptomics, proteomics, or metabolomics) experiment. Pathways overrepresented or enriched in the entities of interest can provide mechanistic insights into the underlying biology of the conditions under study. Many popular pathway databases such as KEGG [1], BioCyc [2], and WikiPathways [3] provide detailed and well-annotated pathways. However, comparisons of pathway databases suggest that no single pathway database is comprehensive [4,5,6]. Further, it has been observed that these databases are partly complementary, and thus it is important for researchers to be able to access pathways from multiple sources simultaneously to gain a more complete picture and not miss possible biological interpretations. The Pathway Architect module in GeneSpring and MPP supports the import and analysis of pathways from these popular pathway databases. In addition, Pathway Architect also supports the import of pathways using standard formats such as BioPAX and GPML.

A lack of standardization in the names and identifiers of biological entities in pathways across multiple pathway databases results in the same entity being cited with different names or identifiers across databases and at times even across pathways within a single pathway database. In some cases, different entities of the same type (gene/protein/metabolite) within a pathway can cite identifiers from different databases as well. Furthermore, in the context of a GeneSpring/MPP experiment, the identifiers associated with entities of interest in the experiment may be different from the identifiers available with the pathway entities. This well-recognized identifier mapping problem poses a major challenge in pathway analysis and limits the matches between the entities from the experiment and their counterparts in pathways.

For example, the metabolite D-glucose (Figure 1) might be known alternatively as dextrose, meritose, or (3R,4S,5S,6R)-6-(hydroxymethyl)oxane-2,3,4,5-tetrol. It has 23 synonyms listed in the Human Metabolite Database (HMDB).
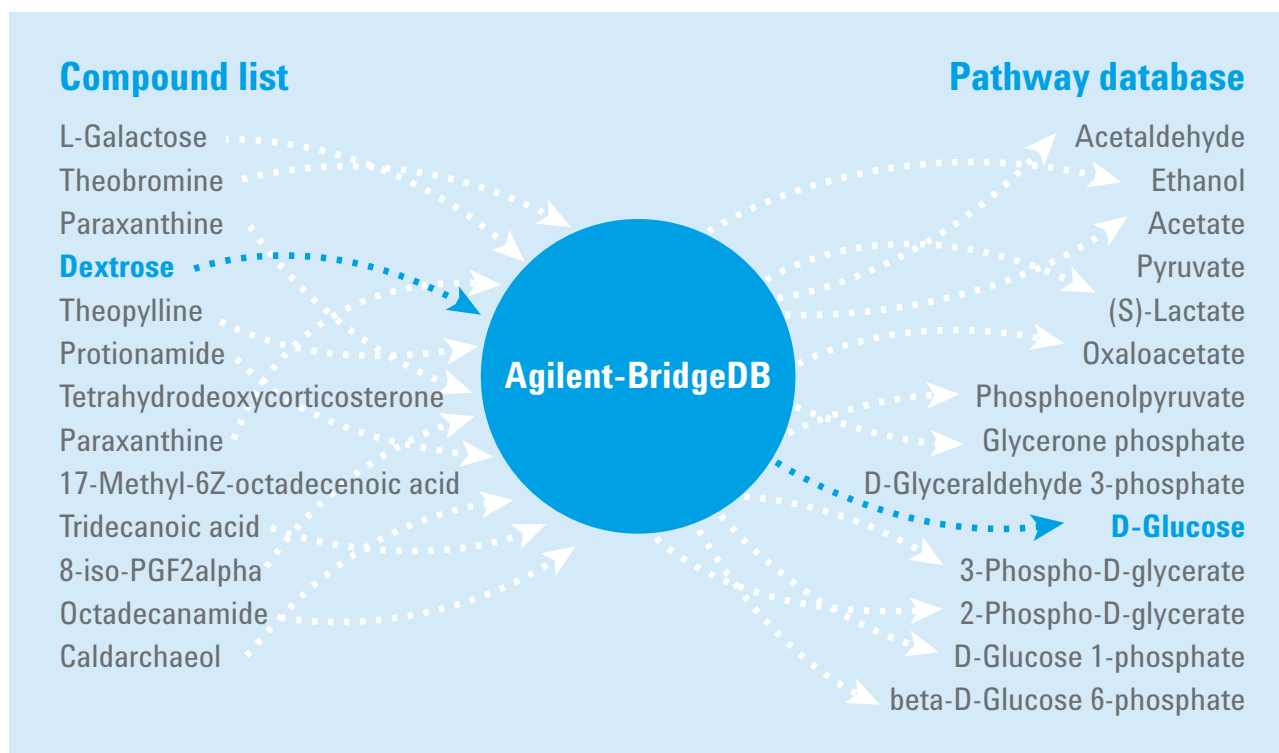


Figure 1. Mapping of metabolites using Agilent-BridgeDB.

To overcome this limitation, GeneSpring/MPP uses a modified version of the BridgeDB software framework [7] to ensure all possible matches between the experiment and the pathway are reported. Mapper files used by the framework provide the mapping between different entity databases to equate an entity in one dataset (pathway) with the same entity in another dataset (experiment). One way to visualize the mapping is in the form of a table where the rows connect all the synonyms and identifiers for an entity (Table 1).

Table 1. D-Glucose aligned with a synonym and some database identifiers.

| Common name | Synonym | KEGG ID | Cas no. | HMDB ID | ChEBI ID |
|---|---|---|---|---|---|
| D-Glucose | Dextrose | C00031 | 50-99-7 | 001222 | 4167 |

There are two types of mapper files currently being used in GeneSpring/MPP-(a) Gene/Protein mapper file and (b) Metabolite mapper file. The Gene/Protein mapper file is organism specific, while the metabolite mapper is common for all organisms. The gene/protein mappers used in GeneSpring/MPP are from the Gladstone Institute and are primarily extracted from Ensembl [8]. The metabolite mapper is developed at Agilent Technologies.

Here we describe four case studies demonstrating the role of Agilent-BridgeDB and the mapper files in enhancing the pathway analysis capabilities of GeneSpring/MPP.



Figure 2. BridgeDB framework in GeneSpring using Agilent metabolite mapper Agilent-BridgeDB and the Gladstone Institute gene/protein mapper to map identifiers across pathways and experiment entities.

## Case Study 1: Mapping different annotations in a pathway and an experiment

Figure 3 shows a pathway in a transcriptomics experiment in GeneSpring. Genes in the experiment are annotated with their Entrez Gene IDs. Table 2 shows an example of the properties available for one of the genes, 'trytophan synthase', in the BioCyc pathway in focus. Genes in this pathway do not cite an Entrez Gene ID, but are annotated with identifiers from other databases. Due to the absence of common identifiers or a bridging mechanism, the pathway does not show any enrichment and the entities do not show any matches with

**A**



**B**



Figure 3. Tryptophan biosynthesis pathway from BioCyc in a transcriptomics experiment. A) without Agilent-BridgeDB and B) with Agilent-BridgeDB. Yellow background color indicates matches with the experiment.

the experiment (Figure 3A). As a result, the pathway is ignored in the analysis.

When the experiment is re-analyzed using Agilent-BridgeDB and the organism specific mapper files, mappings from pathway identifiers to experiment identifiers are retrieved and a match is identified. In the case of tryptophan synthase, the mapping from UniProt/TrEMBL identifier P0A877 (pathway) to Entrez ID 946204 (experiment) is available and is matched in pathway analysis (Figure 3B).

Table 2. Annotations in BioCyc pathway for entity tryptophan synthase.

| Property | Valve | Property | Valve |
|---|---|---|---|
| Cellular location | Cytosol | PDB | 1XCF |
| DIP | DIP-35957N | PR | PRO_000024117 |
| DisProt | DP00252 | PRIDE | P0A877 |
| EcoCyc | TRYPSYN-APROTEIN | PROSITE | PS00167 |
| EcoliWiki | b1260 | Pfam | PF00290 |
| InterPro | IPR013785 | Protein model portal | P0A877 |
| InterPro | IPR011060 | RefSeq | NP_415776 |
| InterPro | IPR018204 | SMR | P0A877 |
| InterPro | IPR002028 | String | 511145.b1260 |
| Label | TrpA | Synonym | Try |
| ModBase | P0A877 | Synonym | TrpA |
| Organism | Escherichia coli K-12 substr. MG1655 | Synonym | Alpha subunit |
| PDB | 1V7Y | Synonym | TSase Alpha |
| PDB | 1WQ5 | Synonym | A protein |
| PDB | 1XC4 | Uniprot/TrEMBL | P0A877 |

## Case Study 2: Mapping of isomers between pathway and experiment

Figure 4 demonstrates a case in which Agilent-BridgeDB enables mapping of specific enantiomers to their D/L form. In this example, both the metabolomics experiment and the metabolites in the KEGG pathway have KEGG Compound identifiers. However, while the experiment contains the KEGG identifier for the D/L form of cysteine (C00736), the pathway cites the isomer specific identifiers: L-cysteine (C00097) and D-cysteine (C00793). Agilent-BridgeDB uses the mappings in the Agilent metabolite mapper to ensure the specific forms of the isomer get mapped to the generic form in the experiment.



Figure 4. Specific isomers in a KEGG pathway are matched with the generic form in the experiment through Agilent-BridgeDB.

## Case Study 3: Mapping multi-omic experiments to multiple pathway databases

Pathways from multiple sources contain complementary information and together are able to provide a more comprehensive picture of biological processes. The ability to map the same entity with different identifiers through Agilent-BridgeDB enables powerful analysis of pathways simultaneously from multiple sources in GeneSpring/MPP. This becomes useful for cases in which pathways from one source cannot be matched with the experiment due to missing annotations. For example, Figure 5 shows the pentose phosphate pathway from two sources, BioCyc and KEGG, enriched in a multi-omics experiment. Metabolites in both pathways could be matched with the experiment. However, proteins in the BioCyc pathway could not be matched with the transcriptomics experiment due to missing annotations, while proteins in the KEGG pathway could be. Thus in the absence of the metabolite mapper files, enrichment of the pentose phosphate experiment from BioCyc would have been overlooked.

**Pentose phosphate pathway**

## BioCyc



**Pentose phosphate pathway**

## KEGG



00030 9/3/13
(c) Kanehisa Laboratories

Figure 5. Multi-omics analysis results for the pentose phosphate pathway from BioCyc and KEGG. Matches with the experiment are indicated by the background color of the entity. Yellow indicates gene/protein matches with the transcriptomics experiment. Blue indicates metabolite matches with the metabolomics experiment.

## Case Study 4: Mapping experiment entities with missing annotations

In some cases, the experiment may have more than one annotation column. It is possible that an entity with a missing identifier in one annotation has been assigned an identifier from another database. For example, Figure 6 shows a genomics experiment with multiple annotation columns: RefSeq Accession, UniGene ID, Ensembl ID, Entrez Gene ID, and Genbank Accession. Note that not all of the database identifiers are present for all entities. Mapping using any single database identifier will invariably lead to loss of matches due to missing annotations. However, pathway analysis in GeneSpring/MPP considers all available annotations for a specific entity in a predetermined order.



Figure 6. Multiple complementary annotation columns in a GeneSpring experiment.

This ensures that entities with sparse annotations are also mapped. In Figure 7 the experiment has the Entrez Gene ID annotation column, but an identifier is not available specifically for the putative 'tubulin' gene. Therefore, Agilent-BridgeDB attempts to match a pathway entity with other available identifiers for this gene. In this case it retrieves a mapping to the UniGene ID and the pathway entity in WikiPathways is matched with the experiment.

| ProbeName | [Untreated] | [Treated] | GeneSymbol | Description | RefSeqAccession | UniGeneID | EnsemblID | EntrezGeneID | GenbankAccession |
|---|---|---|---|---|---|---|---|---|---|
| A_23_P140884 | -0.0043136277 | 0.046764534 | | Homo sapiens, Similar to tubulin, beta, 2, ... | XR_015203 | Hs.513833 | | | BC014971 |
| A_23_P20193 | -0.0043087006 | 0.0032765071 | ARPC1B | Homo sapiens actin related protein 2/3 c... | NM_005720 | Hs.79284 | ENST00000252725 | 10095 | NM_005720 |

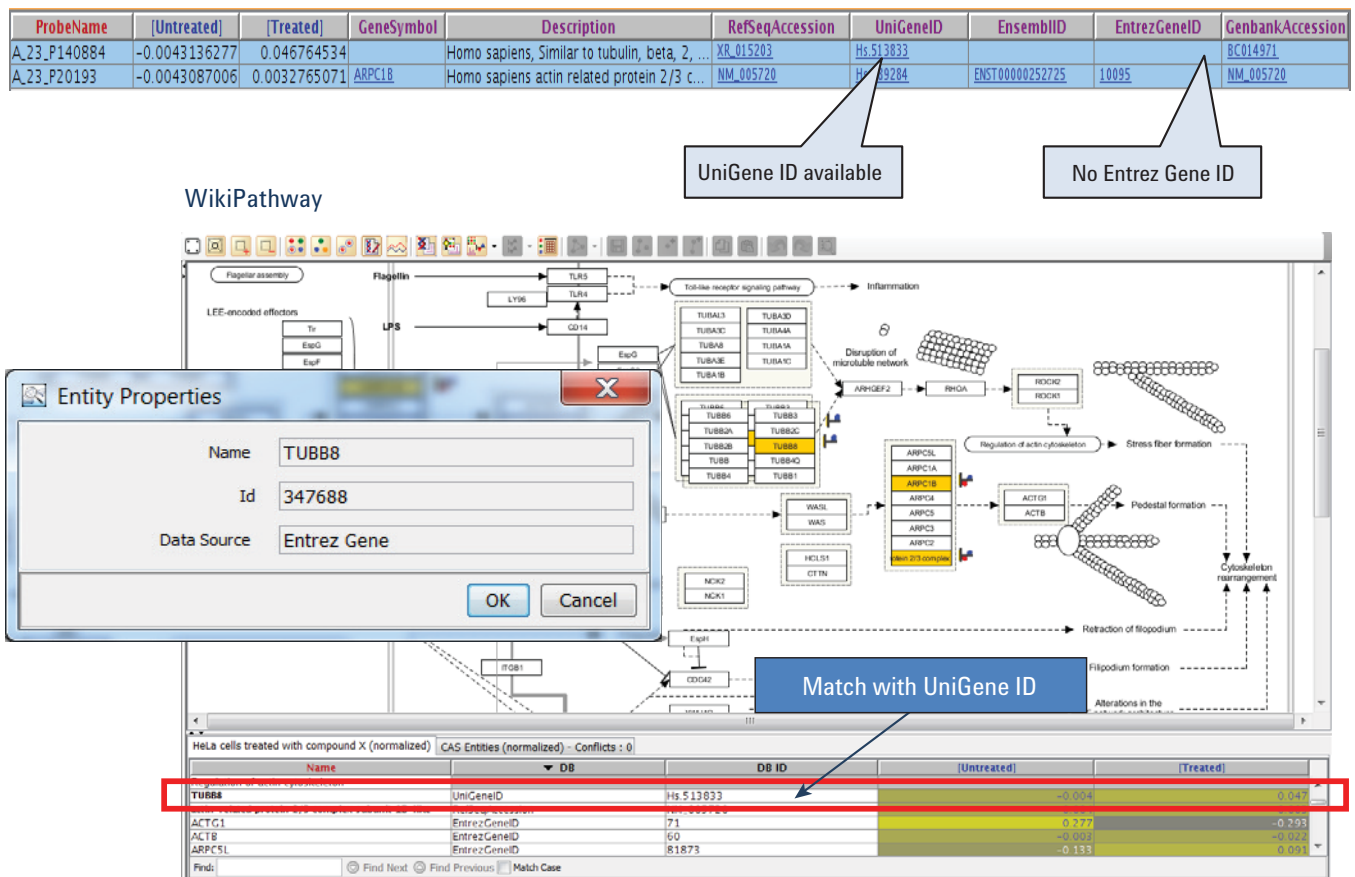UniGene ID available

No Entrez Gene ID

WikiPathway



Figure 7. Pathway entity with Entrez Gene ID is matched to its counterpart in the experiment through its UniGene ID by Agilent-BridgeDB, since the Entrez ID is not available in the experiment.

## Conclusions

Specific examples have been presented across different pathway databases (KEGG, BioCyc, and WikiPathways) and 'omics techniques (genomics, transcriptomics, and metabolomics) available in GeneSpring/MPP. Each of them demonstrated that researchers can get more accurate and comprehensive mappings of their experimental data to pathway databases due to the Agilent-BridgeDB technology. Biological entities that are missing specific annotations in either the experiment or pathway can still be mapped, resulting in more useful information. Multi-omics experiments are more likely to indicate pathways enriched in multiple 'omic technologies since GeneSpring/MPP has mappers for genes, proteins, and metabolites. Successful mapping helps drive research forward by highlighting important pathways and making planning for the next experiment significantly more effective.

## References

1.  Kanehisa, *et al. Nucleic Acids Research*, **42**:D199 (2014).

2.  Caspi, *et al. Nucleic Acids Research*, **38**:D473 (2010).

3.  Kelder, *et al. Nucleic Acids Research*, **40**:D1301 (2012).

4.  Stobbe, *et al. BMC Systems Biology*, **5**:165 (2011).

5.  Soh, *et al. BMC Bioinformatics*, **11**:449 (2010).

6.  Altman, *et al. BMC Bioinformatics*, **14**:112 (2013).

7.  van Iersel, *et al. BMC Bioinformatics*, **11**(1):5 (Jan 4, 2010).

8.  Flicek, *et al. Nucleic Acids Research*, **42**:D749 (2014).

## For More Information

These data represent typical results. For more information on our products and services, visit our Web site at www.agilent.com/chem.

www.agilent.com/chem

**Agilent Technologies**