# CompBio Version 2.5
# User Manual

**Updated July 24 th, 2023**

# Table of Contents

# Section 1 - Account Management

**1.1  Requesting a CompBio Account**

To begin using the CompBio platform, it is necessary to request an account.  All Washington University investigators, researchers, and students may use the platform at no cost.  For research laboratories, we do request that a single account be requested for the lab under the name of the PI.  There is no limit to the number of users within the laboratory that may access the account and, as will be described later in the manual, each user can create their own project folder(s) for organizational purposes.

As the application is web-based, it can be run from any number of browsers and device types.  However, as most testing is done with Google Chrome on PCs and FireFox on Macs we recommend one of these two browsers.

To access the CompBio interface and request and account use the URL - http://at-1.wucon.wustl.edu/

You will then be directed to the login page with the options below:



Click the **Request Account** link and you will be taken to the **Account Registration** page where you will be asked to fill out the pertinent information.

The **Primary Email Address** for the account will likely be the laboratories PI email address as it will be used for general account communications and maintenance in the future.  Once the account request is received by the administration team, an account name with a temporary password will be emailed to that address.

**1.2  Logging in to CompBio**

Once a user account has been created for your lab, you will receive an account name and temporary password.  To login into the system go to the login page using the URL - http://at-1.wucon.wustl.edu/ **Please note that the account id and password are both case sensitive.**

Once you have logged in, you will be taken to the **Project Management** page.  This page contains two main panels.  The first panel titled **Account Administration** to the left of the browser page contains links to all CompBio platform tools as well as your account information.  The second panel contains two to three sub-sections that are specific to the tool you currently have chosen from the first

panel.  The upper-left sub-panel generally contains projects you have completed with that tool and the buttons to the upper-right allow enable the creation of new projects or actions on existing projects. In the case of CompBio, there will be a third lower sub-panel that contains a list of user announcements.  These will typically be notifications of system downtime, updates, and so forth.



## 1.3  Account Maintenance

Clicking the **Account Maintenance** link will take the user to a page that will allow for the **Primary Account Email** address and **Password** to be changed.  This will change either for ALL account users so this should only be done when necessary for the account.

# Section 2 - Project Management

The Project Management component of the CompBio platform is a relatively simple windows-like user interface. While the current functionality is fairly rudimentary, we continue to make improvements and welcome feedback to further enhance the user experience.

## 2.1 Submitting a Project

CompBio is a multi-omics analysis tool and can analyze one or multiple types of entities in a single input list. While CompBio is able to understand several different omics types, there are requirements for the input format on each and, in some cases, limitations on the synonyms that are understood by the platform. This is particularly true for metabolites given the wide range of names they can be identified by in the literature.

### 2.1.1 Entity definitions and limitations

There are five types of entities currently accepted as input for CompBio analysis- Genes, proteins, miRNAs, metabolites, and microbes. All entity lists must be newline delimited. Genes and proteins need to be input as official gene symbols ("ABCA1, IL34"). Microbes need to be input as "genus-species". Without a species name, the microbes will not be included in the analysis. Micro-RNAS need to be in the following format- "mir-1253"as seen in most publications. Metabolites need to be included using their complete names "Carbofuran, Dextromethorphan" etc. A unique feature of CompBio is that the input list can be a mixed list consisting of one or more types of entities.

**Important Note: CompBio can currently analyze a list ranging from <u>5 to 2500 entities</u> in a project. While it is possible to input and run a list larger than 2500, the randomization statistics are currently only computed to this number so larger projects will not have p-values or normalized enrichment scores on the themes.**

In the current version of CompBio, only mammalian gene symbols are considered; specifically human, mouse and rat. Other species are not included yet. However, if the gene symbol used for non-mammalian species in the same as that for mammals, those can get some hits. The metabolite table is extensive but is still in need of some more work because of the nature of those entities. Other entity types, like variants, will be added soon.

### 2.1.2 Expression values

CompBio offers the user an option to include expression values or any other numeric metric corresponding to the entity list. These can be fold changes, p-values or correlation coefficients. These are optional and not used in the creation of the map but are simply carried over into the "theme view". CompBio will generate an identical map even if the expression values are left blank. If the user is input a mixed list, this is a good way to distinguish between gene and protein entries with the same symbol. IL1B as a gene versus IL1B protein. The expression value for IL1B protein can have a "(p)" tag to help distinguish the two.

### 2.1.3 Text file upload

Instead of copy pasting the entity list into the predefined boxes, there is a "file upload" option. A text file (.txt or .csv) containing the entity list and expression values (optional) can be uploaded for a CompBio analysis using the **Choose File** button (1)

Once the entity list is input, user needs to give a meaningful project name (2) and an email address (3) to submit the analysis by clicking **Submit Project**. An email will be sent notifying the user once the CompBio run is complete.

*2.1.4   Batch submission*
Multiple CompBio runs can be submitted at a time using the batch submission option.



Each project needs a name in the Project Name row and entities pasted below. Project test1 is submitted with entities and their expression values. Project test2 and test3 show examples where only entities are input. The progress of the jobs can be tracked using the tracker button on the main page.

## 2.2   Managing Completed Projects

*2.2.1   Project directories*

Once the project is completed and email notification will be sent to the user. A completed project can be moved into a directory/folder within the tool. To create a folder, use the **Create Project Folder** button (4). The new project folder will be visible on the left-hand side panel. To move the project simply select the project and use the **Move to Folder** button (5) to place it into the folder of choice.



### 2.2.2 Annotation markers

Each project map can be annotated by the user. The completed projects which have user defined annotation along with saved profiles will be marked with an orange dot (12) after the project name as shown above.

### 2.2.3 Visualizer optimization

On the project management page, there is a checkbox on the top right side (6). By default, that box is checked. This will limit the visualizer/display of the map to the top 50 themes. Most of the biology generated by the input list is covered within the top 35 or so themes. However, there may be certain cases where a user may want to look for more than 50 themes. This can be achieved by unchecking the box and all themes that can be generated from the input list will be displayed. There are a set of 6 buttons on the top right-hand side of the main panel, below the visualizer optimization check box. These buttons will help the user with basic functionality for the project.

### 2.2.4 Visualize project

Once the CompBio run is completed, the project map can be visualized using the **Visualize Project** button (7). A basic profile generated from the input list will appear in a separate window for the user to annotate and analyze.

### 2.2.5 Delete project

Using the Delete Project button (8) will permanently delete the project. Unless the project was present within the user's account long enough for it to appear in a backup, the project will not be retrievable. Thus, be certain this is acceptable before deleting a project.

### 2.2.6 *Reanalyze project*
The **Reanalyze Project** button (9) will fire a brand new CompBio run using the same input list with the current PubMed and latest version of the PCMM model. The project name will have the current date tagged on the end to distinguish it from the earlier/older run.

### 2.2.7 *Display input list*
The **Display Input List** button (10) will display the input list used to generate the CompBio project.

### 2.2.8 *Announcements*
The **Announcement Button** will display the announcements and updates from the CompBio team.

### 2.2.9 *Move to folder*
The completed CompBio project will be displayed in the main directory. To move to it to the desired folder, the **Move to Folder** button (5) can be used.

### 2.2.10 *Project information*
Once a project is selected, the bottom half of the main panel will display information about the project including the project name, submission date, email used for submission, number of entities.

# Section 3 - CompBio Analysis Projects

Section 3 of this manual will describe biological exploration and knowledge generation from completed CompBio knowledge maps, and the tools required to perform those analyses.

CompBio assembles contextually relevant information derived from an input list of entities (genes, proteins, miRNAs, metabolites, and/or bacterial species). This is accomplished by first identifying concepts enriched in literature that contains the provided entities. The subset of literature that contained the provided entities determines the context for the project. A conditional probability matrix is calculated for each concept versus every other identified concept describing the likelihood of encountering two concepts nearby each other *in the publications returned by the supplied entity list*. Since the probabilities are calculated using only records (publications) that contained the user supplied entities, the biological context is an important factor in the calculation of the distances between concepts. The CompBio algorithm then clusters the closely related concepts into themes.

Below is an overview of the CompBio Workspace. These elements will be discussed in the following sections.
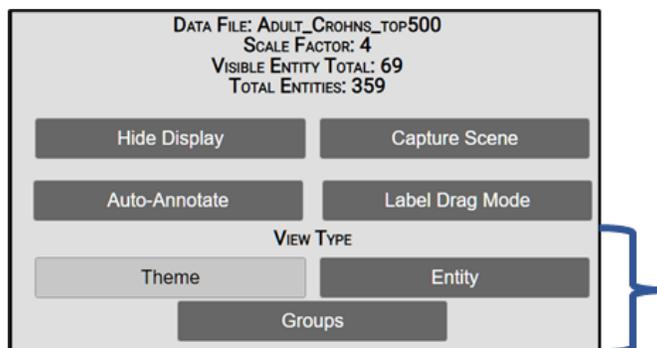


Overview of CompBio Workspace Elements

1. Project panel
2. Control panel
3. Project metrics
   a. Data file
   b. Scale factor
   c. Visible entities
   d. Total entities
4. Hide project panel
5. Auto-annotate project
6. Capture scene
7. Label drag mode on/off
8. Export data
9. Project view selection
10. Theme/Entity detail panel
11. Knowledge map

**3.1      Biological Exploration and Knowledge Generation**

CompBio has been designed to automate what has traditionally been a lengthy analysis process requiring biologists skilled in pathway analysis to wade through lists of differentially regulated entities, assembling groups of genes of similar function to discern pathways and processes of interest.  While current tools can enable this effort to some degree, they tend to suffer from data inflation where a list of hundreds of genes can produce even greater numbers of statistically enriched pathways.  CompBio, on the other hand, assembles enriched concepts into interconnected themes, allowing for rapid interpretation of the knowledge identified and contextually assembled from the input entity list.

*3.1.1      View types – Theme, Entity, and Group*

The CompBio project viewer has three different view types, **Theme**, **Entity**, and **Group**.  Each view type will center the focus of the selection functions in the viewport on that particular type.  The purpose is to allow the user to explore and visualize the data at different levels.   Each of the types and their respective functions will be described in detail in the following sections.  The view types are shown in the upper left section of the **Project Detail** panel and can be selected by clicking on the appropriate box.  **Theme** view is the default view the user will be presented with upon loading a project into the visualizer and will be discussed first. **Entity** view and **Group** view will be discussed at the end of this section (3.1.3 and 3.1.4).



*3.1.2      Theme view*

As described above, CompBio assembles enriched biological concepts into themes based on their associations through a method of natural language processing.  As these themes are assembled without the use of predefined ontologies, the assembled themes generally represent one or more levels of biological process, function, or pathway. In some cases, a theme can contain information from a molecular function to the cell and/or tissue type that function is contextually associated with.

When initially visualizing a project, the Theme view will be displayed by default and the user will be presented with 3-dimensional ball-and-stick model in the central panel that contains the top themes based on Enrichment Score.  Themes are placed close to other themes that they share contextual relationships with.  While the absolute location of a theme is of no consequence, their relative location to each other is very informative.  Proximal themes or theme groups will generally share many entities and describe closely related biology.

*3.1.2.1 Theme annotation*

Once the project is loaded into the visualizer the raw themes need to have annotation applied.  Annotation is a summary description of the concepts contained in the theme.  There are two ways

to accomplish this.  Annotation can be applied through manual annotation, where the user interprets the contents of each theme.  Alternatively, the **Auto-Annotate** function can be implemented which will, through a natural language processing approach utilizing the theme contents and the project context, apply up to three titles to each theme.  Potential titles are sourced from nearly 100,000 biological terms, processes, and pathways.  Each title displays a the confidence score associated with the title autoannotation

**Scores:**

1.2 < Score < 1.5 – Modest confidence (90% accuracy with manual assessment)
1.5 < Score < 2.0 – Marked confidence (97% accuracy with manual assessment)
2.0< Score > 3.0 – High confidence (99% accuracy with manual assessment)
Score > 3.0 – Extremely high confidence (100% accuracy with manual assessment)

An overview of the workflow is shown below.



As can be seen above, three titles have automatically been applied to the theme.  Please note that new projects have the auto-annotation precalculated when the job is first run, but older projects

predating the implementation of this feature will take 1-2 minutes the first time auto-annotation is applied.

This process utilizes the conditional probability matrix described above to not only use theme members to inform annotation, but also closely related terms in other nearby themes. The images below represent graphically the auto-annotation process for a theme related to small-molecule transporters. If one were to only consider the concepts in the theme, a title of "Transporters" may be chosen as a best fit for this theme (shown below in red), but that ignores the contextual information in other nearby themes (shown in green). By allowing contextual inference additional descriptive concepts can be added to the theme title to supplement the understanding of what this theme represents in the context of this experiment. The **theme matrix** shown below represents the conditional probability matrix, showing the probability of encountering two concepts close by each other *in the publications that contained the user-supplied entities*. These associations are specific for the given experiment and are different than the co-occurrence frequency that would be observed in all scientific literature.



This approach allows the annotation process to take the contextual cues for the experiment into consideration while assigning titles to the theme. Given that there are typically several conceptual facets to each theme, the algorithm can assign up to three titles to each theme. In this manner, CompBio produces contextually relevant labels to themes in an automated fashion that can be accepted as-is or modified by the user for group creation.

Manual annotation can also be applied to the theme (theme tools are detailed in 5.4.4) allowing the user to specify a custom name for any or all themes.

### 3.1.2.2 *Concept-entity views*

The entities provided by the user when the CompBio project was created identify concepts that are clustered into themes. The concepts within a theme, and the entities that identified them, can be visualized by clicking on the theme of interest which brings up a **Theme Members Window**, shown below.

13

**Theme Members Window**



The theme concepts are listed on the left-hand side of the **Theme Members Window**, and the entities are listed on the right. The Centrality scores and Enrichment scores are discussed in sections 3.1.8 and 3.2.1, respectively.

*3.1.2.3  Source literature*

In CompBio, entities identify literature (abstracts and full text articles) and enriched concepts within those documents. The links between entities, concepts and the source literature is preserved in CompBio, allowing the user to drill down to the underlying literature that contains the entity (or one of its synonyms) and the concept by clicking on the concepts or entities within the **Theme Members Window**. This workflow is depicted below after clicking on a theme.

[**Section continues on next page with large image**]

14

**Literature Scanning Page**

Clicking on any entity listed on the right-hand side of the **Theme Members Window**, shown above, will spawn a **Scope Box** as shown below.



**Scope Box**

There are two options here, and this will determine the scope of the entity contribution scores returned.

- **Entity-To-Theme**: This option returns a list of concepts *from within the selected theme* and the score describing how much the selected entity contributed to the identification of that concept.
- **Entity-To-All**: This option returns a list of concepts *from the entire project* and the score describing how much the selected entity contributed to the identification of that concept.

The difference here is the scope of the results returned. This choice about scope does not apply when clicking a concept in the Theme Members Window because concepts only occur uniquely in a single theme, whereas entities, as mentioned previously, can contribute to multiple themes. Typically, users are trying to understand how a given entity contributed to the selected theme, so the Entity-To-Theme option is more typically chosen. The Entity-To-All option informs the user on the role that entity plays in the entire project.

The literature results are returned in a **Literature Scanning Page** that contains condensed and formatted text. This condensed text lacks punctuation and is formatted to describe entity and concept occurrence in the knowledge map.

Below is a key describing the formatting contained in the **Literature Scanning Page:**

- **Red**: any occurrence of the entity (or any entity synonym) selected in the **Theme Members Window**
- **Blue**: any occurrence of the concept selected in the **Entity Window**.
- **Purple**: concepts appearing in the literature results that belong to the selected theme.
- **Green**: concepts contained in the current project that appear in the literature results, but do not belong to the selected theme.

The **Literature Scanning Page (**next page) summarizes the literature results from the selected entity-concept pair, but also serves as a gateway to the original underlying literature. The workflow to drill down into the original publications is shown below and simply requires clicking on the **See Original Document** link for the given abstract.

**ChatGPT summaries**: These summaries are generated by ChatGPT based on the information contained in the top abstracts identified by CompBio. They are meant to provide a short, rapid assessment of key information associated with the selected entity and concept in the context of the given theme. The number of abstracts can vary but typically falls into the range of the 8 to 10, if that many are available from the analysis. Please note, that the query restricts ChatGPT to simple summarization of the provided text to reduce the probability of ChatGPT introducing false statements. Within these constraints falsehoods have been rare, though not zero, so it is recommended to review the abstracts for confirmation if important information is being gleaned from the summary.

## ChatGpt Generated Summary

**Selected Entity: Red**    **Selected Concept: Blue**    **Concepts in Theme: Purple**    **Concepts in Project: Green**

**Entity: ppara**    **Concept: pparalpha**

**Chatgpt Summary of Top Abstracts:**

The text discusses the role of the gene **ppara** and the biological concept of **peroxisome proliferator-activated** receptor alpha (PPARalpha). **PPARalpha** is a nuclear receptor that regulates lipid and glucose metabolism and is activated by **peroxisome** proliferators, a group of chemicals that induce predictable responses including liver tumors in rodents. The text highlights that of the three isoforms of **PPAR** (alpha, beta, and gamma), **PPARalpha** is essential for the responses induced by **peroxisome** proliferators. Activation of **PPARalpha** leads to the transcriptional activation of genes involved in **fatty** acid **oxidation** and **cytochrome P450** enzymes. **Peroxisome** proliferation and increased oxidative stress in liver cells are also observed. The text also mentions that disruption of the gene **ppara** in mice leads to spontaneous **peroxisome** proliferation and liver tumors. Additionally, the text discusses the interaction between **PPARalpha** and other proteins, such as HSP90 and XAP2, and their role in regulating **PPARalpha** activity. Overall, the text emphasizes the importance of **PPARalpha** in regulating lipid metabolism and its potential role in liver tumor development.

# Literature Scanning Page



**Source Literature**

*3.1.2.4 Concept map*

A theme represents a clustering of the concepts that were positioned in three-dimensional space based on the conditional probability of those concepts co-occurring with each other in the literature identified by the entities. As such, the concepts within a theme therefore can be represented with a **Concept Map** describing their level of cooccurrence with each other. Concepts are ranked in the **Theme Members Window** according to how central they are to the theme (discussed in detail in 3.1.2.5.1). The **Concept Map** helps to understand the central concepts in the theme enabling the best manual theme annotation possible. The workflow for accessing the **Concept Map** is shown below. Note that the underlying literature is also accessible by clicking on edges in a workflow similar to that described in section 3.1.2.6 of this manual.



*3.1.2.5 Image search*

To inform the user about processes and pathways associated with the theme, a modified Google Image search has been implemented to return (primarily) scientific results related to the concepts contained in the theme. This can be useful in manual theme annotation as well as obtaining images for presentation purposes, or to obtain review articles related to theme concepts. This image searching utilizes all of the concepts within the theme, but customized searches are also possible (described in 3.1.2.5.2). The **Image Search** workflow is shown in the following image.

**Theme Members Window**

**Image Search**

### 3.1.2.5.1 *Centrality*

The main CompBio knowledge map is created by arranging the project concepts in three-dimensional space based on their contextual cooccurrence, and then clustering them into themes. This approach means that each theme contains a subgraph of the concepts in 3D space. Accessing the Concept Map is described in 3.1.2.4. A **Centrality** score is calculated describing the distance of the concepts from the center of the theme, which is used to rank the concepts. Entities with a low distance from the center (low **Centrality** score) can imply high relevance or could describe a midpoint between two similar but distinct processes clustered into a single theme. In the **Concept Map** shown in 3.1.2.4, "fatty" is the most central concept, but the map shows "fatty" as a mid-point between two distinct but related concept branches, namely PPAR/Peroxisome concepts on the left and lipid metabolism on the right. It can be seen in 3.1.2.1 that the auto-annotation labels assigned to this theme correctly capture both of these biological facets of the same theme.

### 3.1.2.5.2 *Custom image search*

In addition to the **Image Search** described in 3.1.2.5, it is also possible to perform a custom search from within the **Concept Map** described in 3.1.2.4. This allows the image search to be performed on a subset of the concepts and can be helpful in themes that have more than one biological facet. Another powerful feature of this tool is that **Additional Query Term(s)** can be specified in a comma delimited list. This allows the user to understand the selected themes within the context of the additional terms. An example search is shown below; the search is initiated by clicking **Begin Search**.

### 3.1.2.6 Edges

The knowledge map consists of entities and concepts clustered into themes represented by spheres. These themes can be, and often are, interconnected with each other by edges. An edge can potentially be formed so long as there is at least one gene that contributes to one or more concepts with an enrichment greater than 2.0 in both themes. Greater numbers of shared genes will cause the given edge to increase in thickness and low numbers of shared genes will similarly cause a thinner line to be displayed.

Selection of an edge in the project window will spawn an **Edge Entity Window** that contains the theme labels and the entities shared between the two themes. As described in 3.1.2.6, edges can potentially be drawn if the same entity contributes to concepts in two themes, so long as the distance threshold criteria (described in 5.4.3) are met. These shared entities are displayed in the **Edge Entity Window** as shown below**.**



*3.1.3.6.2 Source literature*

The information that drove concept enrichment for this project and enabled edge formation is traceable back to the original abstracts. This enables the biologist to understand and drill down into the publications that identify a given concept being related to the selected gene. The figure below demonstrates clicking on an edge to display the shared contributing genes in the **Edge Entity Window** and drilling down to the abstracts that link that gene to the selected concept displayed in the **Edge Concept Window.** Clicking on any entry within the **Concept** will immediately take the user to the Abstract Scanning page with the ability to access the original abstracts or publications. See section 3.1.2.3 for more detail.

**Edge Entity Window**

Adherens Junctions/Catenin Signaling(4.2)
Cadherins/Cell-Cell Junction(4.2)
Connexins/Cx43(2.9)
- & -
Intestinal Brush Border/Enterocyte(5.1)
Apical/Basolateral Exchangers(3.2)
43 Shared Entities

| Entity | Enrichment Scores | |
|--------|------|------|
| slc15a1 | 101.7 | 108.7 |
| slc36a1 | 83.7 | 54.4 |
| ush1c | 48.6 | 58.8 |
| tjp3 | 449.5 | 30.9 |
| slc5a1 | 26.1 | 192.9 |
| slc9a3 | 25.1 | 383.7 |
| erbin | 31.2 | 22.9 |
| cybrd1 | 31.7 | 207.9 |

Entity-To-Themes    Entity-To-All    Close ⊗

**Selected Entity:** slc15a1

| Concept | Score | Theme |
|---------|-------|-------|
| (p)caco-2 | 1253.61 | Adherens Junctions/Catenin Signaling(4.2) Cadherins/Cell-Cell Junction(4.2) Connexins/Cx43(2.9) |
| caco-2 | 1253.61 | Adherens Junctions/Catenin Signaling(4.2) Cadherins/Cell-Cell Junction(4.2) Connexins/Cx43(2.9) |
| (p)enterocytes | 353.53 | Intestinal Brush Border/Enterocyte(5.1) Apical/Basolateral Exchangers(3.2) |
| enterocytes | 353.53 | Intestinal Brush Border/Enterocyte(5.1) Apical/Basolateral Exchangers(3.2) |
| transepithelial | 347.38 | Adherens Junctions/Catenin Signaling(4.2) Cadherins/Cell-Cell Junction(4.2) Connexins/Cx43(2.9) |
| enterocyte | 259.92 | Intestinal Brush Border/Enterocyte(5.1) Apical/Basolateral Exchangers(3.2) |
| basolateral | 234.47 | Intestinal Brush Border/Enterocyte(5.1) Apical/Basolateral Exchangers(3.2) |
| apical | 207.89 | Intestinal Brush Border/Enterocyte(5.1) Apical/Basolateral Exchangers(3.2) |
| brush | 194.48 | Intestinal Brush Border/Enterocyte(5.1) Apical/Basolateral Exchangers(3.2) |
| border | 92.75 | Intestinal Brush Border/Enterocyte(5.1) Apical/Basolateral Exchangers(3.2) |
| permeability | 77.87 | Adherens Junctions/Catenin Signaling(4.2) Cadherins/Cell-Cell Junction(4.2) Connexins/Cx43(2.9) |
| (p)microvilli | 39.22 | Intestinal Brush Border/Enterocyte(5.1) Apical/Basolateral Exchangers(3.2) |

**Edge Concept Window**

### 3.1.3   Entity view

The **Entity View** allows the user to interact with the knowledge map from the standpoint of the entities, as opposed to the themes and concepts (**Theme View,** discussed in 3.1.2). The **Entity View** allows the user to see which entities contributed to concept identification across the themes, and to what degree. In the workflow shown below, after the **Entity** button is clicked, an entity

can be selected from the **Theme/Entity Detail Panel**, and the visible themes will be haloed with colors representing the **Enrichment Score** of the selected entity for those themes. The color scale key is displayed in the **Project Panel** when the **Entity** button is clicked. This approach is useful to understand the selected entity contribution to all visible themes.



### 3.1.4 Group view

Once the theme annotation process is complete, either by auto-annotation or manual annotation (5.1.4), themes can be bundled into groups of similar biological function (5.4.5). The **Groups** button then allows the user to interact with the knowledge map at the group level. The workflow shown below details entering **Group View** and selecting a group. In this mode, group-level coloring can be applied (5.4.5). Additionally, newly selected themes can be added to the group using the **+ Add to Group** button, removed from the theme by clicking the "x" bullet preceding the theme name, or the entire group can be deleted using the **Delete Group** button. Note that the **Delete Group** button does not delete the raw themes, but rather removes the group name from the **Group View** and discards any group-level formatting.

## 3.2 Understanding the Data

The focus of CompBio is to identify the enriched concepts given a list of input entities. To enable this objective, statistics and metrics are employed to allow the biologist tasked with data interpretation to assess the relative significance and enrichment levels of the identified concepts, as well as deciding objectively when to conclude the annotation process.

### 3.2.1 Significance

When a project is run in CompBio the project is compared against thousands of randomized comparator projects of similar entity list size. This allows the algorithm to calculate enrichment scores for concepts and p-values for the themes as compared to random noise. The level of precision reported in these metrics is determined by the number of randomized projects that have been run. In the current version of CompBio (v 2.5), p-values are reported down to the .001 level. Any p-value lower or equal to 0.001 are simply reported as "<0.001".

In the **Theme Detail Panel** on the lower left side of the workspace, significance and enrichment metrics are displayed as well as concepts contained within the themes and auto-annotation (if applied), as shown below.

Regulation of lipid metabolism by Peroxisome proliferator-
activated receptor alpha (PPARalpha)(4.7)
Synthesis of very long-chain fatty acyl-CoAs(4.4)
PPAR/RXR/Beta-Oxidation(4.0)

Enrichment Score: 80562.15
Entity Total: 47
PValue: <0.001
NEScore: 2.521
Concept Total: 12

fatty,peroxisome,hepatic,acyl-coa,oxidation,proliferator-
activated,long-
chain,metabolic,carnitine,ppar,pparalpha,beta-
oxidation,peroxisomal,steatosis,palmitoyltransferase,diet,pr
oliferator,dietary,fasting,medium-chain,
(p)peroxisomes,peroxisomes,triacylglycerol,lipolysis,polyun
saturated,rxr,droplets,
(p)lipid_droplets,adiponectin,acylcarnitine,palmitoyl-

There are five key metrics displayed for each theme that enable the biologist to determine the robustness of the theme. Weaker scoring yet significant themes tend to provide a smaller degree of useful information towards the interpretation of the dataset as compared with the higher scoring themes. This will be discussed more thoroughly in the next section describing entity mapping efficiency.

- **Enrichment Score:** magnitude of enrichment of all of the concepts in the given theme over randomly generated themes of the same rank with similar numbers of input entities
- **Entity Total:** number of entities that contributed to identifying concepts in the theme
- **P-Value:** p-statistic calculated to a describe the confidence that the given theme is significantly enriched over random noise. So long as the **NEScore is >= 1.3**, p-values as high as 0.1 should be given consideration as a 0.05 cutoff tends to be overly strict with suitably enriched CompBio themes
- **NEScore:** the given theme's **Enrichment Score** normalized to the enrichment of the same theme rank in projects created with random entity lists within the same size range
- **Concept Total:** total number of concepts contained within the given theme

*3.2.2    Entity mapping efficiency*

One of the strengths of CompBio is the efficiency with which it is able to map the provided entities to biological concepts and themes. Other tools in widespread use tend to produce large amounts of statistically relevant pathway results that need to be reviewed and analyzed to determine what is meaningful. CompBio addresses these issues by producing a graph of the percentage of user-supplied entities mapped to concepts, which can be accessed through the **%Entities Mapped** button on the bottom right-hand side of the workspace (graph shown below).

As can be seen in the above graph, themes are listed across the X-axis and the cumulative number of entities mapped is shown on the Y-axis. By theme 23, over 90% of the provided entity list has been mapped to at least one biological concept. This provides a logical and meaningful cutoff when deciding what data to consider and what to set aside. This is important in the annotation workflow, as it is possible for a project to have over 50 significant themes and limiting the interpretation to just the primary mappings can save significant amounts of time in the annotation process.

### 3.2.3 General quality control

It is important to understand the level of biological signal in a dataset at the outset of analysis, and CompBio has several tools to enable this. The **Theme Detail Panel** contains metrics that describe the levels of enrichment of the biological signal as well as the level of significance. The **Enrichment Score** describes the degree of signal enrichment of the concepts in a given theme, over random noise. This can be thought of as a "fold change" over random noise. The theme order is defined by **Enrichment Score.**

The **NEScore** is based on the **Enrichment Score** but has been normalized to randomly generated themes of the same rank number. Therefore, the **NEScore** for theme 1 in a given project is calculated by normalizing the **Enrichment Score** for that theme against thousands of theme 1 **Enrichment Scores** from knowledge maps generated from randomly created lists of entities. This puts the themes on a common scale, regardless of rank, allowing for comparison of relative enrichment.

The **P-Value** is a significance calculated from the normally distributed enrichments of concepts identified in randomly created lists of entities.

After testing with both random entity lists, positive control pathways, and well-characterized published datasets, the following guidelines are suggested for theme evaluation:

- **NEScore** >= 1.3
- **P-Value** < 0.1

If a theme fails to meet both of the above criteria, it is suggested to view the theme as failing to meet general quality control thresholds.

### 3.2.4 *Visuals and display control*

The **Control Panel** located in the upper right hand of the workspace contains the tools for modifying the knowledge map. In the current version of CompBio, there are 14 primary sections in the **Control Panel**, shown below. Primary sections 2-9 are expandable, with subfunctions dropping down accordion-style after the user clicks on the primary section, which is denoted with a small triangle proceeding the primary section title.

| | |
|---|---|
| ① | Profile     Default ⌄ |
| ② | Load Profile |
| ③ | ▸ Profile Tools   ❓ |
| ④ | ▸ Scale Factor   ❓ |
| ⑤ | ▸ Filters   ❓ |
| ⑥ | ▸ Theme Tools   ❓ |
| ⑦ | ▸ Group Tools   ❓ |
| ⑧ | ▸ Visual Setting   ❓ |
| ⑨ | ▸ Export Data   ❓ |
| ⑩ | Set Center Point |
| ⑪ | Reset Center Point |
| ⑫ | Capture Scene |
| ⑬ | ❓ Help |
| ⑭ | Close Controls |

A brief description of each of the tabs is provided below.

**1. Profile:** CompBio makes use of profiles to allow multiple views of the same knowledge map to be generated. This will be described in more detail in the **Profile Tools** section. The **Profile** pull-down box allows for the selection of a previously created profile or allows the user to revert to the original **Default** view.

**2. Load Profile:** Applies the profile selected.

**3. Profile Tools:** The profile tools primary section contains a drop-down with a single option. This option allows the user to either create a new profile or to delete the current profile. If a user profile is loaded, another primary section appears on the **Control Panel** allowing the user to **Save Progress**.

**4. Scale Factor:** This allows the user to modify the repulsive force that the themes apply on each other, modifying overall separation of the graph. This is useful if the knowledge map is tightly packed and greater readability is desired.

**5. Filters:** Controls for limiting the themes displayed in the knowledge map based on various criteria (described in detail in section 4.1.1.1)

**6. Theme Tools:** Controls for modifying the name and text in a theme as well as specifying theme-specific edge display (described in detail in section 5.4.4)

**7. Group Tools:** Controls for defining, naming, coloring, and display of groups of themes (described in detail in section 5.4.5)

**8. Visual Setting:** Text display tools, background, and various other visual control options (described in detail in section 5.4.6)

**9. Export Data:** Tools for compiling numerical data that describes the entire knowledge map, or a subset of the map (described in detail in section 5.4.7)

**10.**        **Set Center Point:** Allows the user to modify the center point around which the knowledge map rotates

**11.**        **Reset Center Point:** Reverts the center point back to the default location.

**12.**        **Capture Scene:** (described in detail in section 5.4.7)

**13.**        **Help:** Hyperlink to this document

**14.**        **Close Controls:** Hides the **Control Panel**

# Section 4 - CompBio Project Comparison

It is often desirable to compare two or more data sets to gain biological understanding of similarities and/or differences in them. Different human disease cohorts, animal models and human disease, different treatments within a model or disease, or time points within a treatment course are just a few examples of the conditions for which these comparisons are often sought.  Within the CompBio platform, projects can be compared in multiple ways.  Simple comparisons can be done at the concept or entity level of any two projects, and much more sophisticated contextual comparisons can be completed on libraries of projects with the Assertion Engine tool.  This section will describe the utility and directions for each type of comparison.

## 4.1  Basic Project Comparison

### 4.1.1    Concept comparison

CompBio projects can be compared at the concept level by first creating a concept filter and then applying that filter to a given project.  While these filters are often utilized with one or both projects from which they are created, they can be used with any CompBio project(s) within the users account.  As these filters are created from a simple list comparison of concepts or entities from the comparator projects, the results should be viewed with care.  Even though a given concept may exist within the top themes of two different projects, the context, or other concepts within those themes, may be different across the two projects.  Thus, the user is primarily looking for themes in which numerous concepts are shared between the projects.   Alternatively, entity level comparison will create a list of shared entities between the two projects and will create a new CompBio project from that shared list.  The user can then observe the differences and similarities between the child project and the two parent projects.

### 4.1.1.1  Filter creation

The first step in project comparison is to create a filter that can be applied.   To enable a project comparison, select the **Project & Filter** link under the **Comparison** tools section in the **Account Administration** panel.  This will direct the user to the **Compare Projects** panel.  To create the project level filter, click the **Concept Level (New Comparison Filter)** button.



The **Overlap Project** section will now be visible in the lower sub-panel.  Enter the name of the filter you wish to create in the **Project Name**.  You can then use the two drop down lists under **Base Project** and **Compare To Project** to select the two projects for comparison.  Once selected, click the **Create Project** button.  This will create an intersection of the concept lists that will be saved under the project name provided and can be used as a concept filter for any completed CompBio project.

Clicking the **Display Overlap List** will show the complete list of intersecting concepts from the two projects. This is the list that will be applied as filter in the visualizer. Selecting a filter from the **Existing Comparison Lists** will display the basic information about that list including the creation date, the base and overlap project and total concepts in the overlap.

*4.1.1.2 Filter application*

Concept filters will now be visible in **Filters** tab of the **Controls** panel on the far right of the visualizer screen of any loaded CompBio project. To access the filter, open the Filters tab and click on the **Overlap Filter** white box. This will bring up a display of all filters available to apply to the project. Select the appropriate filter from that list.



Once the filter is selected, the user will then need to decide in which mode to utilize it. Just below the **Overlap Filter** option is the **Filter Mode** option. *Intersection* mode will display concepts that are shared with the filter list, and *Difference* mode will display concepts that are not shared with the filter list within the visualizer screen. **Note that these two modes are mutually exclusive.** Once the mode is selected, click the **Apply Filter** tab. Themes will now redisplay with the concepts that match the selected filter mode.

In the example above, Intersection mode was used with the filter list created from a NASH cohort and NAFLD cohort. The concepts still visible within the themes intersect with the filter list applied. The colored halos that are now present around the themes indicate how many of the concepts originally in the theme overlap with the filter list in heat-mapped fashion. Just below the theme annotation label are a pair of numbers. The first indicates the number of concepts in that theme that match with the filter lister and the second is total number of concepts in the theme. As the filter list was created by comparing NASH and NAFLD cohorts, this view indicates that themes associated with FXR/Bile Acid, Sterols/LXR, LDL/Atherogenesis, and Lanosterol Synthase activity are strongly preserved between NASH and NAFLD. Additionally, themes associated with Villi, NASH, and P450/Xenobiotics are modestly preserved, and themes related to Prostaglandins/Eicosanoids and NADP are weakly preserved. It should be noted that if a theme does not have a single preserved concept with the filter, the entire theme will disappear from the view. This was true in the case above as theme related to Flavin-containing MonoOxygenase activity was unique to NASH and no concepts overlapped with NAFLD. If the given filter had been applied in Difference mode, the view would have flipped. All currently displayed concepts would disappear, and the absent concepts would appear. Application of the filter in this manor would present the concepts that are unique to the NASH when compared to NAFLD project.

### 4.1.2    Gene comparison and new project creation

While the New Comparison Filter function will simply compare the two full lists of concepts from the completed projects and allow the user to view that filter applied as described above, the **Gene Level (New Comparison Project)** utility will create an intersection of the entities that were used to create those two projects and will generate a third, unique CompBio project using that intersected list. When comparing biology from different species, such as human and mouse, or from two different technology platforms where the entity lists may not be identical, it will generally be more enlightening to use the Comparison Filter. However, when different conditions are applied to common background, such as different drug treatments or a time course, it may be more desirable to compare at the entity level. Thus, the platform enables both types of comparisons.

As with the concept level comparison, the entity overlap project generation can be found under the **Project & Filter** link under the **Comparison** tools section in the **Account Administration** panel.  After clicking the **Gene Level (New Comparison Project)** button, the **Create An Overlap Project** dialogue will appear in the lower sub-panel.  Utilization is identical to that of the concept level comparison described in 4.1.1.1.



Clicking the **Create Project** button will take the user to the CompBio **Project Submission** page with the **Entity List** pre-loaded with the supplied **Project Name** and entity intersection list.  To run the project, enter the Email notification address and click the Submit Project button.  The user will be notified when the new project is complete.



## 4.2  Assertion Engine – ML Based Contextual Comparison

The Assertion Engine program within the CompBio platform enables the user to perform a comparison of projects, in pairs, to assess what biological concepts and processes are conserved between them.  Unlike the **Project & Filter** functions that simply create an intersection of identical concepts or entities between two project that can then be applied as filter within the CompBio project visualizer, the Assertion Engine uses a machine learning-based, pattern similarity identification utility to detect patterns of contextually preserved biology.  This means that the engine does not simply look to see if a concept exists within both projects, but pattern of connectivity between that concept and other concepts is preserved.  This assessment is performed across the entirety of the concept space yielding a global map, with local regions, of preserved biological connectivity.  Conceptually, this can be visualized in the example(s) of how similar or dissimilar different representations trees, or components of trees, are in the following figure.  In this case, even though components are highly overlapping, the context of some components change.  The Assertion Engine can recognize both the similarities and dissimilarities without ever training to understand what a tree or its components are.  It merely identifies the level of preservation across components in any comparison and can, statistically, determine if the level preservation is significant or not.

Complex Knowledge Pattern Comparison
Simplistic Example of Assertion Example

From a biological perspective, this ability allows the Assertion Engine to identify preserved biology in a unprecedented fashion across projects. Examples of how the engine has been used to date include comparison of different human disease cohorts with a common diagnosis, comparison of human disease to animal models, comparison of different intervention treatments in human or models, comparison of stages of progression in human disease or models, and time points within a time series. Though the engine can only compare two projects at a time, libraries of project can be created and passed to the engine to compute a full block of comparisons across many conditions or related projects of interest.

### 4.2.1    *Creating projects and libraries*

Accessing the Assertion Engine is done by clicking the **Assertion Generation** link within the **Comparison** section of the **Project Management** page. The inputs to the Assertion Engine are existing CompBio projects or project libraries. To begin an Assertion Engine run using existing completed projects click the **New Assertion 2.0** button in the upper right section of the **Assertion Generation** page.



This will open the dialogue setting up the engine run and generating a comparison.

Select the base and secondary projects to compare using the two drop down lists. These lists will show the names of all completed projects available from within the user's account. Enter the name of the Assertion project. This can be any name that does not conflict with an existing Assertion Engine run name. Finally, click the **Generate Assertion** button. A typical comparison will take a minute or less.

If the user wishes to compare a large number of projects to a project or another group of projects, this can be done by first creating one or more libraries from existing CompBio projects. To create a library, first click the **Show Libraries** button in the upper right section of the **Assertion Generation** page. This will bring up a new list of available functions in the upper right section of the page while showing any existing libraries in the upper left section.



To create a new library, click on the **New Assertion Library** button. This will initiate new dialog in the lower portion of the page.



First, enter the new libraries name in the **Library Name** dialogue box. Then, using the **Library Projects** drop down list, begin to select the project to be added to the library. It should be noted that while there is currently no limit on the number of projects that can be added to a library, once number of projects starts to exceed 8 to 10 projects, the resulting heatmap matrix can become a

35

bit difficulty to read.  Therefore, we recommend keeping the number of projects in a library to 10 or less.  As the user adds project by selecting them from the drop-down list and then clicking the **Add To Library** button, they will appear in the lower right section under the Assertion Library List header.  If the user decides to remove one or more projects, simply click the **red X** next to the projects name before clicking the **Create Library** button.  Once the list is complete, click the **Create Library** button to generate the library.  It will now appear in the **Assertion Eligible Libraries** list and can be selected in the **Generate New Assertion 2.0** dialogue.   While individual project assertion comparisons generally take less than one minute to run, large library comparisons can take several minutes to an hour depending on the total number of assertions that result.  Once the assertion run is completed the resulting library or project level assertion will display in the **Completed Assertion Generation(s):** list.  Select the desired assertion within that list and click the **Visualize Assertion** button to view the results.

*4.2.2    Global similarity assessment*

Upon clicking the **Visualize Assertion** button, the user will be taken to the **Global Similarity** comparison page showing project level comparison in a heatmap and dendrogram format.  The pseudo heatmap section also displays the overlap value and a p-value that evaluates the significance of that overlap based on the size of the input projects.



The example in the above figure illustrates the comparison between 4 windows of genes selected from a list that was identified as differential (down regulated) between adult Crohn's Disease subjects and non-IBD controls.  The Top 500 list contains the 500 genes with p-value < 0.01. As these genes were the most significantly different from the non-IBD controls, they are expected to represent base biological truth for the data set.  The Bottom 500 list are genes with 0.01 < p-value < 0.05. While these genes are still nominally significant and would be expected to contain some relevant biology, the level of noise within the list is anticipated to be higher than the Top 500 gene list.  The non-significant 500 (Nonsig500) gene list are the next 500 genes with 0.05 < p-value.  These genes may contain some signal but likely contain more noise than signal.  The Random 500 are genes randomly selected that are not contained with any of the first three lists. Based on these criteria, the first three lists would be expected to contain some degree of similar biology to one another but with increasing levels of noise.  The Random 500 list would be

expected to have no significant biological overlap with any of the first three. As can be seen in the figure, this is exactly what the Assertion Engine identified.

Each cell in the pseudo-heatmap contains a global overlap or similarity score for the project pair and a p-value indicating the level of significance for that overlap. Clicking on the **Overlap Score** button from the comparison of the top 500 ACD genes to the bottom 500 ACD genes will bring up the biological similarity map for those two projects. The thickness of the lines between concepts is an indicator of the level of contextual preservation between those concepts and what they, in turn, associate with. The colored regions are rough boundaries for sets of concepts that are tightly conserved. These will tend to be concepts associated with a particular process, pathway, cell type or other organized biological construct.

| Dendrogram rows | | Dendrogram columns | | Dendrogram all | |
|---|---|---|---|---|---|
| | | **ACD_Top_500** | | **ACD_Bottom500** | |
| ACD_Top_500 | | 1.0000, | | 0.1024 , | P < .001 |
| ACD_Bottom500 | | 0.1024 , | P < .001 | | |
| ACD_Nonsig500 | | 0.0880 , | P < .001 | 0.072 , | P < .001 |
| ACD_Random500 | | 0.0330 , | P > .9 | 0.0472 , | P = .3 |



Each row of the pseudo-heatmap is based on the same comparator project. As such, the structure and regions of the map will remain the same as the user clicks the different **Overlap** button within the different cells in that row. This is not true for the columns as different cells in a column change the aspect to a different project.

As, in this case, the library was compared to itself, there is a visible green diagonal where each project was compared to itself. For these cases, the maximum possible overlap value of 1.0 is observed and, as such, no p-value is displayed as it is recognized as a self-comparison.

The **P-value** button also provides a graph to demonstrate just how far from the randomized range a given overlap is. The calculation of the p-value is done empirically by comparing randomized data sets of different sizes to real-world data sets of size of the base project in each case.

## Section 5 – Visualizer Controls and Navigation

### 5.1 Visualizer navigation

#### 5.1.1 *Workspace elements overview*
Once the project is complete,the "visualize project" button (see section 2.2.3) will display the CompBio map. The visualization shows 3 panels- Left hand side panel gives the theme details with scores, p-values and concepts. The right-hand side shows the control panel for visual setting options and the middle panel is the actual 3D map generated from the input list. Refer section 3.1 for a detailed description.

#### 5.1.2 *Themes, concepts, and entities*
The panel on the left-hand side shows the themes ordered by the total enrichment score from theme 1 to theme 50 (default settings is 50 themes. Can be opened up to greater than 50 using the visualizer checkbox (6)-see Visualizer Optimization section 2.2.3.)
Each theme will display the concepts that make up the theme. If the theme is clicked, (Theme2 in the figure below, it shows an orange hue), a separate box will appear in the main map (2) which shows details with concepts and genes that brought those concepts. If expression data or any numeric value is provided as input, those values will be displayed here next to the enrichment scores for each entity.

### 5.1.3    3D controls: pan, zoom, rotate

The central map can be maneuvered with the mouse- pan, zoom and rotation operations can be performed to get the perfect view for image capture. Left mouse click allows you to rotate; right mouse click for pan and the mouse wheel can be used to zoom in and out.



### 5.1.4    Auto-annotate

One of the recent features within CompBio is the ability to auto-annotate themes. Using the **Auto-Annotate** button (1) themes can be assigned a name based on the conditional probability matrix using the concepts within the theme as well as related themes. Details of the auto-annotation are described on section 3.1.6.

Many themes can have more than one label assigned. Using the **First Line Only** checkbox (3) only the top scoring label can be displayed. The labels can be moved and aligned for best display using the **Label Drag** button (2).

### 5.1.5    Capture scene

Once the themes are annotated and the final visual display of the map is created, the image can be saved using the **Capture Scene** button (4,5). The button is available within the control panel. Once the button is clicked a separate window with the map and a **Download Button** (1) will appear. A .svg image file will be saved on your desktop than can be edited using any graphic software such as Adobe illustrator.

### 5.1.6    Label drag mode

The labels assigned by auto-annotation can be moved around and aligned for best display using the **Label Drag** button (2).

### 5.1.7    Export unmapped entity data

Once the input list is submitted, the user can find out valid and mapped entities using the **Entity Mapping Data** button (1)



Once this button is clicked, a.csv file will be downloaded on the computer showing the details of the input list as shown. The list will have 5 columns. First is the input list. Second column shows invalid entities- those entities that did not have a core memory in the model, i.e. those genes that are not recognized. The third column shows valid entities- those that the tool recognizes and can use in the analysis. The fourth column shows the mapped entities- those that are valid and meet a threshold minimum 3 abstracts. the last column shows the unmapped entities- those that are valid but do not meet the analysis threshold.

unmapped_entities_1664932081462

| Input List (97) | Invalid (13) | Valid (84) | Mapped (77) | Unmapped (7) |
|---|---|---|---|---|
| ac058822.1 | ac058822.1 | acta2 | igfbp3 | isoc2 |
| ac099489.1 | ac099489.1 | aif1 | igfbp5 | linc01833 |
| acta2 | h2bc18 | ak1 | grb10 | nme3 |
| aif1 | h2bc6 | ankrd1 | tgfb2 | senp3-eif4a1 |
| ak1 | h2bp2 | arfgef3 | sstr2 | tceal5 |
| ankrd1 | h3-2 | arhgap23 | slc2a1 | zfp28 |
| arfgef3 | h3c1 | arhgap4 | dlk1 | znf585b |
| arhgap23 | h4c1 | atp5md | edn1 | |
| arhgap4 | h4c2 | atp5mg | mt-atp6 | |
| atp5md | h4c6 | bahcc1 | mt-co3 | |
| atp5mg | rn7sl4p | cdh6 | mt-nd6 | |
| bahcc1 | rnu2-2p | cdkn2b | mt-nd3 | |
| cdh6 | vax1 | celf2 | mt-nd1 | |
| cdkn2b | | col8a1 | mt-nd2 | |
| celf2 | | cox6a1 | mt-cyb | |
| col8a1 | | crip2 | mt-co2 | |

5.1.8    Filter creation - See section 4.1.1
5.1.9    Filter application - See section 4.1.2

## 5.2   Assertion Engine Figures - See section 4.2

## 5.3   Data exports for Third Party Tools
The CompBio tool allows the user to export theme mapping data that be easily formatted into supporting material for manuscripts. These export tables can be downloaded from the export data pull down in the control panel.

*5.3.1 Export selected (publication table)*
Generate a csv file with the filename given in the **Data file name** box. This file will have entities and their scores for all the themes that are selected from the map.

| Entity ID | Entity Score | Expression Value | Theme Name (Score) |
|---|---|---|---|
| (e)zfp36 | 4708.360235 | -6.38 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)zfp36l1 | 2136.744062 | -1.56 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)auh | 367.561846 | -1.51 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)dcp1a | 155.707367 | -1.79 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)ybx1 | 104.616552 | -1.43 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)nxf1 | 84.758610 | -1.39 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)eif3a | 56.418078 | -1.63 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)ythdf1 | 47.796287 | -1.71 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)eif4e | 46.581231 | -1.35 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)eif3i | 45.306897 | -1.72 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)dhx9 | 38.308198 | -1.64 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)eif3d | 33.464730 | -1.71 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)rps3 | 32.364950 | -1.34 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)rbm3 | 30.675109 | -1.55 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)hnrnpc | 29.877187 | -1.68 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |
| (e)arid5a | 27.480190 | -1.3 | Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) (99762.51) |

*5.3.2 Export selected (raw data)*
Exports theme details for all the themes that are selected within the map. Each column corresponds to a theme. For each theme the column will show all the metrics associated with that theme. The raw and normalized enrichment scores, p-value, total counts for entities and concepts and individual entity and concept that make up that theme.

| Name=TIM23 mitochondrial import inner membrane translocase complex( 4.9) | Name=Tristetraprolin (TTP, ZFP36) binds and destabilizes mRNA( 7.0) |
|---|---|
| Score=207974.49 | Score=99762.51 |
| Coordinates=X: -8.07,Y: 14.03,Z: -4.00 | Coordinates=X: -8.31,Y:-23.23,Z:-38.05 |
| Entity_Total=70 | Entity_Total=38 |
| Concept_Total=21 | Concept_Total=4 |
| PValue=<0.001 | PValue=<0.001 |
| NEScore=3.823 | NEScore=2.590 |
| Entity=(e)timm23(1867.791927) | Entity=(e)zfp36(4708.360235) |
| Entity=(e)timm17a(1298.066706) | Entity=(e)zfp36l1(2136.744062) |
| Entity=(e)tomm40(1242.639155) | Entity=(e)auh(367.561846) |
| Entity=(e)tomm7(679.099883) | Entity=(e)dcp1a(155.707367) |
| Entity=(e)immt(404.891225) | Entity=(e)ybx1(104.616552) |
| Entity=(e)ndufb8(179.786080) | Entity=(e)nxf1(84.758610) |
| Entity=(e)ndufb6(153.115144) | Entity=(e)eif3a(56.418078) |
| Entity=(e)ndufs8(148.153272) | Entity=(e)ythdf1(47.796287) |
| Entity=(e)ndufs7(145.476518) | Entity=(e)eif4e(46.581231) |
| Entity=(e)higd2a(109.646909) | Entity=(e)eif3i(45.306897) |

Export all raw will have the same format except that it will have 50 columns, one for each of the 50 themes. All these tables can be used to import into R or other data processing tools and used as supplementary material for manuscripts.

*5.3.3 Export bar graph*
A bar graph can be created using the **Bar graph** option under the **Create Figure** selection from the control panel.
There are several options presented to the user. Either a bar graph of all themes can be created or a selected few. The themes can be selected from the map just by clicking the theme spheres. The other way is to give a range of p-values or NEScores to plot all these passing those cutoffs. An image as shown in the panel will appear with red bars showing the NEScores on the x-axis and

the theme labels on the y axis. The p-values will also be displayed. The svg file containing the graph will be downloaded on the desktop. These can be incorporated into any manuscript figures.



## 5.4    Visualizer Fine Controls
While many of the concepts being discussed in this section have been touched on previously, section 5.6 will delve into greater detail around use of the tools contained within the **Control Panel**.  A figure detailing Control Panel components and a brief description of their functions is shown in section 3.2.4 of this manual.

### 5.4.1    Creating, naming, and saving a user project profile
When a project is initially visualized, the **Default Profile** is displayed that contains up to 50 themes, unless the Optimization Mode option is deselected on the **Project Management** page, at which point hundreds of themes could potentially be displayed.  The **Default Profile** represents the starting point for your project, the raw, clustered output from CompBio.  As the annotation process, described in section 3, proceeds, it is necessary to save the progress made.  This is accomplished in CompBio through the use of profiles.  Having multiple profiles allows multiple users to work on the same knowledge map and to create multiple views of the same project.  To create a new profile, click the **Profile Tools** tab in the **Control Panel.**  A drop-down tab opens with the **Create New** option.  Three new drop-down tabs appear.  In the **New Name** tab enter a suitable name for the new project profile and click **Set New Name**.  A popup in the middle of the workspace will appear saying "Progress Saved", indicating that the profile has been successfully saved.  This profile is now associated with the project and in the future can be selected from the

**Profile** pull-down and clicking **Load Profile**.  Similarly, a profile can be deleted using the **Delete Profile**.

### 5.4.2  Visualizer scaling

The **Scale Factor** tab allows the user to modify the repulsive force that the themes apply on each other, modifying overall separation of the graph.  This is useful if the knowledge map is tightly packed and greater readability is desired.  To increase the separation of the graph, enter a value greater than the default value of 4, and click **Initiate Scaling**.

### 5.4.3  Visualizer filters

Filters that adjust the themes displayed are located in the **Filters** tab in the **Control Panel.** Clicking the **Filters** tab will open the filters as shown below.  Tabs 7 – 10 are filters for project comparison and are described in detail in sections 5.1.1-2.



1. The **P-Value** filter limits the themes displayed based on p-value, which was described in section 3.2.1.  Dragging the slider to the left increases the stringency, requiring lower p-values for a theme to remain on the screen.  Any changes to any of the filters is automatically applied to the project, but it is recommended to click **Save Progress** under frequently to ensure changes made to the knowledge map are save under the current profile.
2. The **Score(%)** filter limits the themes displayed based on their enrichment values.  Dragging the slider to the left increases the stringency, requiring higher enrichment scores to remain on-screen.
3. The **NEScore** filter limits the themes displayed based on their normalized enrichment score.  Dragging the slider to the left increases the stringency, requiring higher normalized enrichment scores to remain on-screen.
4. The **Distance (%)** filter limits the edges displayed based on their intra-theme distance in three-dimensional space.  Dragging the slider to the left increases the stringency, requiring smaller distances between themes for an edge to be shown.
5. The **Depth (Z)** filter limits the themes displayed based on their position in three-dimensional space.  Themes farther from the viewer will be removed from view as the slider is dragged to the left.
6. The **Depth On/Off** check box toggles the **Depth (Z)** filter on and off.

7. The **Overlap Filter** pull-down box allows for the selection of previously created comparison filters to limit the themes displayed based on concept overlap. See section 5.1.1-2 for details.
8. The **Filter Mode** pull-down box allows for the selection of **Intersection** or **Difference** options showing themes that overlap with concepts from the overlap filter or are unique form concepts in the filter, respectively.
9. The **Filter Halos** check box toggles theme halos that highlight the degree of overlapping concepts between the current project and the filter. The color scale is, from most relative overlap to least: red, orange, green, and blue.
10. The **Apply Filter** button must be pressed to apply overlap filter pull-down options.

*5.4.4    Theme tools*

During the annotation process, theme attributes are modified using controls in the **Theme Tools** tab in the **Control Panel.** Below is a figure marking the seven controls that reside within this tab.



1. The **Theme Name** text box allows the user to specify a custom name for the selected theme.
2. The **Set Name** button applies the specified **Theme Name**. Any changes to the **Theme Name** will only be saved to the profile after clicking **Save Progress**.
3. **Hide Names** toggles the theme names displayed on the knowledge map on and off.
4. **Hide Concepts** toggles the theme concepts displayed on the knowledge map on and off.
5. The **Color** text box allows the user to either specify the theme color in hex color code or to pick the color using the mouse from a color palette that pop up on mouse hover.
6. The **Set Color** button applies the color chosen in the **Color** text box. Any changes to the theme color will only be saved to the profile after clicking **Save Progress**.

7. The **Always Show** button causes an edge to be drawn between two (and only two) selected themes, so long as they share at least one common gene, regardless of the maximum distance threshold specified with the **Distance (%)** filter (see section 5.6.3).

### 5.4.5    *Theme group creation and naming*
During the annotation process themes are bundled into similarly colored groups of themes with similar biological function using the controls in the **Group Tools** tab in the **Control Panel.** Below is a figure marking the six controls that reside within this tab.



1. The **Group Name** text box allows the user to specify a custom group name for the selected themes.
2. The **Set Group** button applies the specified **Group Name**.  Any changes to the **Group Name** will only be saved to the profile after clicking **Save Progress**.
3. **Hide Group** hides the currently selected group on the knowledge map.
4. **Hide Concepts** toggles the theme concepts displayed on the knowledge map on and off.
5. The **Color** text box allows the user to either specify the theme color in hex color code or to pick the color using the mouse from a color palette that pop up on mouse hover.
6. The **Set Color** button applies the color chosen in the **Color** text box.  Any changes to the theme color will only be saved to the profile after clicking **Save Progress**.

### 5.4.6    *General visual settings*

The **Visual Settings** tab contains tools to modify various visual attributes of the knowledge map not contained in the themes or groups tabs. Below is a figure marking the ten controls that reside within this tab.



1. The **Theme Text** slider allows the user to specify a custom text size for the themes in the knowledge map.
2. The **Concept Text** slider allows the user to specify a custom text size for the concepts in the themes in the knowledge map.
3. **Label Text** allows the user to specify a custom label text color for the theme labels in the knowledge map.
4. **Label Background** allows the user to specify a custom label background color for the theme labels in the knowledge map.
5. The **Label Score** checkbox allows the user to toggle the numerical scores shown in the theme labels on and off.
6. The **Bold Labels** checkbox applies bold text formatting to the theme labels in the knowledge map.
7. The **First Line Only** checkbox restricts the auto-annotation theme labels, which can range from zero to three lines (section 3.1.6) to only a single line of annotation.
8. **Wireframe** checkbox converts the textured themes to wireframe rendering. This can be useful if the effect is desired for aesthetic reasons, or to improve performance on low-performing computer hardware.
9. The **Background** pulldown allows the user to select one of several preset backgrounds for the knowledge map.

10. The **Background Color** text box allows the user to either specify the background color in hex color code or to pick color the using the mouse from a color palette that pop up on mouse hover.

*5.4.7    Data exports*

The **Export Data** section of the **Control Panel** contains tools to export numerical and graphical data from the knowledge map for further analysis.  Below is a figure marking the ten controls that reside within this tab.



1.   The **Export Selected (Publication Table)** button allows the user to export theme-level data detailing the entities and enrichment scores for selected themes in the knowledge map.
2.   The **Export Selected (Raw Data)** button allows the user to export theme-level data detailing the entities, concepts, and enrichment scores in a less formatted configuration for selected themes in the knowledge map.
3.   The **Export All (Raw Data)** button allows the user to export theme-level data detailing the entities and enrichment scores for all themes in the knowledge map.
4.   The **Data File Name** text box allows the user to specify the exported file name.
5.   The **Export Graph** button allows the user to export a bar graph of the selected theme enrichment scores.

6.  The **Entity Mapping Data** button allows the user to export the original entity input list, as well as invalid genes (genes that had insufficient literature to contribute to concept identification in the knowledge map), mapped genes, and unmapped genes.
7.  The **Set Center Point** sets the current visual center of the knowledge map as the point of rotation for the graph.
8.  **Reset Center Point** button reverts the point of rotation of the knowledge map back to the original settings.
9.  The **Capture Scene** button allows the user to capture the knowledge map in a new window as an SVG image with tools to move and edit the graph to suit the user's needs. Exportable as a publication-ready SVG image.
10. The **Help** button is a link to this document.

### 5.5    Publication Examples

Below are several examples of how published or submitted manuscripts have incorporated CompBio results to convey findings in publication-quality figures.

1.      Simple table detailing the major CompBio findings.

## Table 1

CompBio in silico analysis of bulk RNA sequencing data showing significantly altered pathways in livers from NAMPT OE and WT controls.

|  | SIRT-dependent | SIRT1-independent |
|---|---|---|
| Increased | PI3K/AKT signaling | Neutrophil Chemotaxis |
| NAMPT OE vs. WT | NFκB Signaling | Complement Pathway |
|  | T cell Regulation |  |
| Decreased | NASH | De novo lipogenesis |
| NAMPT OE vs. WT | Lipid metabolism | Hepatic Steatosis |
|  |  | Retinoic acid |
|  |  | Adipogenesis |

**Nat Commun. 2022; 13: 1074**

2.  A heatmap-style summary of themes across experiments.

| Enriched Themes | BRAFmut versus | | EGFRmut versus | | IDH1mut versus | |
|---|---|---|---|---|---|---|
| | EGFRmut | IDH1mut | BRAFmut | IDH1mut | BRAFmut | EGFRmut |
| **Immune Function** | | | | | | |
| Histocompatibility | red | red | blue | | blue | |
| Interleukin signaling | red | red | blue | | | |
| T-cell modulation | red | | blue | red | | blue |
| Allergic-type inflammation | red | | blue | | | |
| Chemokine signaling | | | | red | | blue |
| Mucin activity | red | | blue | | | |
| **Trophoblast-Related Pathways** | | | | | | |
| Trohoblast-like phenotype | red | red | blue | | blue | |
| Trohoectoderm development | red | | blue | | | |
| **Tissue Remodeling** | | | | | | |
| Matrix metalloproteinase activity | red | red | blue | blue | | red |
| Cellular junctions | | red | | | blue | |
| Collagen remodeling | | red | | | blue | |
| Angiogenesis | | | | red | | blue |
| Aggrecan remodeling | | | | red | | blue |
| Integrin expression | | | | red | | blue |
| Mesenchymal phenotype | | blue | | | red | red |
| **Protein Processing** | | | | | | |
| Endoplasmic reticulum / protein folding | red | red | blue | | blue | |
| Golgi trafficking | red | red | blue | | blue | |
| Protein glycosylation | red | red | blue | | blue | |
| Ubiquitin mediated decay | | blue | | | red | |
| **Cell Division / Renewal** | | | | | | |
| Cell cycle / mitosis | blue | blue | red | blue | red | red |
| Regulation of apoptosis | blue | red | red | | | |
| Pluripotency | red | | blue | blue | | red |
| **Signalling Cascades** | | | | | | |
| MET signaling | red | red | blue | | blue | |
| MAPK/ERK signaling | blue | blue | red | | red | blue |
| SMAD / TGF-B signaling | blue | | red | | | |
| WNT signaling | | blue | | blue | red | blue |
| ERBB2 signaling | | | | red | | blue |
| **Gene Expression Regulation** | | | | | | |
| Epigenetic modification - methylation | red | | blue | | | |
| RNA mediated silencing | blue | | red | | | |
| mRNA splicing | blue | blue | red | | red | |
| **Neuron Funciton** | | | | | | |
| Neurdegeneration | blue | | red | | | |
| Neurotransmission | | blue | | | red | |
| Neurogenesis | | blue | | blue | red | red |
| Synapse function | | | | blue | | red |

**Sci Rep, 2021 Oct 8;11**

3. Drilling down from themes (in wireframe) down to genes and concepts.



**EBioMedicine, 2021 May;67:103347**

4. Detailed pathway analysis and knowledge maps



**Fig. 3.** Plasma proteomic biological themes correlated with changes in CAZyme gene abundances after consumption of the fiber-snack prototypes. (*A*) Schematic illustrating the distribution of plasma protein projections along singular vector 1 (SV1) from a CC-SVD analysis of the abundance of CAZyme genes versus changes in levels of plasma proteins after fiber snack consumption. Proteins with SV1 projections along the tails ($\alpha = 0.1$) of the distribution are highlighted (green and red); these proteins, belonging to the 10th and 90th percentile groups, were analyzed using CompBio to identify biological themes enriched in each of the two percentile groups for each treatment. (*B*) Diagram illustrating pathways in which microbiome CAZyme-correlated plasma proteins (boldfaced) are involved in three interrelated biological themes (TGF-β/BMP-mediated fibrosis, P38/MAPK-associated immune biomarkers, and VEGF-mediated angiogenesis). (*C, D*) Network of biological themes identified from the CC-SVD and CompBio analysis. Themes are depicted as spheres; the number in each sphere corresponds to the theme enriched after pea or orange fiber treatment that is listed in *SI Appendix*, Dataset S6D. The size of a sphere is proportional to the CompBio enrichment score of its theme. The thickness of the lines connecting themes is proportional to the number of proteins shared between them. Purple spheres represent themes related to TGF-β-BMP signaling (fibrosis), p38/MAPK immune biomarkers, and VEGF-mediated angiogenesis that were enriched in the plasma proteomes of pea and orange fiber study participants (see Table 1).

**Proc Natl Acad Sci, 2022 May 17;119(20)**

5. Drilling down from knowledge maps to themes across experiments, to expression of individual genes from themes.



Figure 3

**Silver et al., Manuscript in preparation**

**a**

PCA2 21.6%

PCA1 29.0%

PCA3 14.2%

**b**

WT
ank/ank

M M F F F F F M F M

**c**

21
19.1
17.2
15.3
13.4
11.5
9.6
7.7
5.8
3.9
2

*ank/ank* Avg (log2)

GLUL
MYC
SLC11A2
DUSP6
TNFSF15
UGT1A10
VPS35
LEPR
GALNT11
MAN2B1
CCK
EID1 AUH
SLC35D2

UP
DOWN

2  3.9  5.8  7.7  9.6  11.5  13.4  15.3  17.2  19.1  21
WT Avg (log2)

**d**  Downregulated themes in *Ank^ank/ank* NP

Synostosis of carpal bones (11 genes)
Ichthyosis skin (12 genes)
COL17A1-related bullous pemphigoid associated (7 genes)
Defective SLC24A4 causes hypomineralized amelogenesis imperfecta (AI) (11 genes)
Simpson-Golabi-Behmel syndrome (10 genes)
Paraxial mesoderm development (16 genes)
Eph/ephrin signaling related to synostosis (8 genes)
Spinocerebellar ataxia type 1 (18 genes)
tRNA (m1A) methyltransferase complex (14 genes)
LPA/Lysophospholipid Signaling (11 genes)
3-Methylglutaconic aciduria (8 genes)
Lipofuscinosis/Lysosomal Degradation (4 genes)
Pyrophosphate-dependent phosphofructokinase complex (16 genes)
bHLH/DNA Binding (13 genes)
BMAL1/Clock Circadian Regulation (12 genes)
Transport of Ribonucleoproteins into the Host Nucleus
Integrin-mediated Adhesion/Signaling (25 genes)
Oxygen-dependent proline hydroxylation of HIF-1A (24 genes)
Filamin binding (24 genes)
Retromer, tubulation complex (24 genes)
Galactosylation of collagen propeptide hydroxylysines by procollagen galactosyltransferases 1, 2. (11 genes)
PP-InsP5 kinase activity

**e** Retromer, tubulation complex

lyst, cpe, man2b1, acp2, rabgap1l, tmem106b, atg12, cln3, rps9bp, kif13a, stx5a, stxbp6, snx9, snx27, vps35

CompBio Entity Score

**f** BMAL1/Clock Circadian Regulation

amelx, mef2c, eid1, slc12a2, ppp1r3c, tbx6, ppip5k2, id3, bhlhe41, dbp, nr1d1, per3

CompBio Entity Score

**g** Pyrophosphate-dependent phosphofructokinase complex

ppp1r3b, fubp1, ppp1r3c, eef1g, cd36, p4ha2, epas1, cat, jmjd8, acadvl, vps35, egln3, me1, hk2, pfkl

CompBio Entity Score

**h** Defective SLC24A4 causes hypomineralized amelogenesis imperfecta

fmod, fgfr3, kcnk2, adm, fgf11, galnt3, sparc, xylt1, slc39a13, col1a1, amelx

CompBio Entity Score

**i**

| GLOBAL OVERLAP (AE) Senmayo Panel vs ANK | | | | | | NP_DOWN 0.0621, P = .02 | AF_UP 0.0435, P = .4 | AF_DOWN 0.0826, P < .001 |
|---|---|---|---|---|---|---|---|---|
| Theme (SenMayo Panel) | Score | Entities | Concepts | <0.001 | NEScore | NP_DOWN | AF_UP | AF_DOWN |
| Insulin-like growth factor (IGF) activity regulation by insulin-like growth factor binding proteins (IGFBPs) | 202917 | 39 | 11 | <0.001 | 6.475 | 0% | 0.00% | 9.09% |
| Tissue Inhibitor Of Metalloproteinases (TIMP) associated ECM remodeling | 119024 | 75 | 25 | <0.001 | 4.67 | 48% | 4.00% | 12.00% |
| urokinase-type plasminogen activator proteolytic cleavage product | 104754 | 70 | 22 | <0.001 | 4.646 | 9% | 0.00% | 4.55% |
| VEGF and VEGFR signaling network | 87660.9 | 67 | 25 | <0.001 | 4.271 | 12% | 0.00% | 24.00% |
| Lymphotoxin-A related TNF and receptor activities | 84824 | 49 | 11 | <0.001 | 4.461 | 9% | 0.00% | 0.00% |
| Protein kinase A (PKA) and RPS6KA5 (MSK1) phosphorylates p65 (RELA) subunit | 76225.8 | 104 | 40 | <0.001 | 4.282 | 3% | 2.50% | 5.00% |
| Increased nuchal translucency | 74102 | 27 | 19 | <0.001 | 4.42 | 11% | 0.00% | 5.26% |
| prostaglandin E2 receptor EP2 subtype | 69471 | 47 | 19 | <0.001 | 4.382 | 0% | 5.26% | 15.79% |
| Tumor necrosis factor (TNF) pathway | 68291.9 | 91 | 38 | <0.001 | 4.533 | 3% | 7.89% | 0.00% |
| granulocyte macrophage colony-stimulating factor receptor activity | 66268.5 | 78 | 40 | <0.001 | 4.615 | 3% | 0.00% | 7.50% |
| GP130/IL6 Signaling | 63395.4 | 56 | 4 | <0.001 | 4.615 | 0% | 100.00% | 75.00% |
| IL-7 signaling pathway | 60894.9 | 43 | 4 | <0.001 | 4.627 | 50% | 0.00% | 0.00% |
| CD95 death-inducing signaling complex | 58216.1 | 36 | 12 | <0.001 | 4.601 | 17% | 16.67% | 50.00% |
| JNK/MAPK/ERK Pathway | 56594 | 69 | 38 | <0.001 | 4.651 | 3% | 13.16% | 15.79% |
| chemokine (C-X-C motif) ligand 12 production | 55859.6 | 66 | 11 | <0.001 | 4.765 | 0% | 0.00% | 9.09% |

## 5.6 Methods Description for Manuscripts

**Platform Overview**
CompBio (v2.5) is an artificial intelligence platform that enables the interpretation of omics and multi-omics data performed through a systematic analysis of literature pertaining to a list of biological entities (genes, proteins, miRNAs, or metabolites) provided the by user. Entities can be obtained from differential expression, correlation, or any number of relevant analysis methods. The first major component of the platform extracts and stores knowledge from PubMed abstracts and full-text articles using contextual language processing that is not restricted to fixed pathway and ontology knowledge bases. The extracted knowledge is maintained in a Persistent Contextual Memory Model (PCMM) for future use. The second major component is the generalized signal detection engine that can aggregate, normalize and compute across data from any stored dimension of the PCMM. Within the engine, biological concepts that are enriched with the input entities are identified and aggregated into higher-level biological themes (pathways/processes, cell types and structures, etc.) that are themselves inter-connected through their identified associations. The resulting knowledge maps are interactive with fully traceable information.

**Scoring and Significance**
Conditional probability analysis forms the foundation of gene, concept, and theme-level enrichment scores. The enrichment scores, in turn, are normalized and assessed for significance empirically, utilizing large, randomized comparator groups similar in construct to the query group. As such, the reported Normalized Enrichment Scores (NES) represent the magnitude to which the concepts/themes are enriched above random, and an empirically derived p-value identifies the likelihood of achieving that NES by chance.

**Cross-Project Contextual Analysis – Assertion Engine v1.0**
Fully contextualized analysis of multiple CompBio projects can be generated using the Assertion Engine (AE v1.0) component of the platform. The Assertion Engine is a machine-learning based engine capable of identifying preserved patterns between biological concepts of the compared projects and their inter-concept relationships. In other words, if a concept is present in both projects the AE determines how similar or different that concept's relationships are with the other concepts in the respective data sets. Strongly preserved groups of concept/inter-relationship associations will form high scoring sub-maps within the overall analysis provided to the user. A global score, representing the complete contextual biological similarity between the two is also computed. A score of 0.0 represents no similarity and a score of 1.0 represents complete similarity. P-values are also computed for the significance of the global similarity score.

**Project and Theme Annotation**
Within CompBio, biological themes and their relationships are computed in a fully automated fashion that is independent of existing ontologies. The biology described by the concepts in the themes can often be obvious to the user, however, it is generally desirable to have high-level annotation, at a process or pathway level, for the themes that is commonly used and recognized by a broad audience. To accommodate this, the auto-annotation system utilizes a language processing method that leverages the theme contents and over 100,000 biological terms, processes, and pathways to produce optimal annotations in a fully automated fashion. However, unlike traditional pathway/process tools, the CompBio annotation system not only uses the contents of a given theme when annotating it, but also considers the contents of neighboring themes as well. This provides overall annotation within the project that is, generally, more accurate and contextually specific. The same principles are utilized to generate project-level annotation as well.