

IMPC annotation pipeline

Annotation pipeline in the International Mouse Phenotyping Consortium (IMPC) is an existing data assignment project with a goal to associate phenotypic observations to the genetic modification. Here we explain the steps that are taken to select the best Mammalian Phenotype (MP) term to the genetic modification in mice when a significant difference (typically at the level of 0.0001) from the baselines observed.

Annotation pipeline and the analysis framework

The IMPC annotation pipeline (IMPC-AP) assigns MP terms to the significant genetic effect. The genetic effect at the IMPC is specified by three statistical analysis platforms that are designed in the IMPC statistical pipeline through OpenStats software. Here we break the annotation pipeline by the type of the input data and the analysis frameworks.

Annotation table

The annotation pipeline in the IMPC requires a reference table that summarizes the available terms for an IMPC parameter. This can be retrieved from IIMPreSS however to remove the dependency to the live servers, the IMPC-AP utilised an offline version of the file called *Annotation Indexer* in this document. This file is available from XXXXX.

Continuous data – Linear mixed model

Continuous data such as tail length, tibia length etc. in IMPC is analysed by linear mixed model, implemented in the software package OpenStats. The continuous measurements are more informative than the other types in that aspect that the direction of change can be determined by the effect size. Here we summarised the steps to assign MP terms to the continuous measurements.

| From the statistical results | From the Annotation Indexer |
|---|---|
| <ol style="list-style-type: none">Overall effect (both sexes)<ul style="list-style-type: none">if $p\text{value} \geq \text{threshold}$ → Assign no MP termif $p\text{value} < \text{threshold}$ (II)<ol style="list-style-type: none">If effect size > 0 → Increase termif effect size < 0 → Decrease termif effect size $= 0$ → Steady termsimilar steps in 1 apply for Male effect (III)similar steps in 1 apply for Female effect (IV) | <p>Filter for</p> <ol style="list-style-type: none">Pipeline_stable_idProcedure_groupParameter_stable_id <p>Get available MP terms (I)</p> |
| Find matches between I and II, III, IV | |
| Notes | <ul style="list-style-type: none">If increase or decrease effect detected then ignore ABNORMAL MP term.Generally accepted threshold by the IMPC consortium is 0.0001 |

Continuous data – Reference Range plus

Due to the complexity of the data not all continuous data can be analysed by linear mixed model. Alternatively, there are many cases in the IMPC that are analysed by Reference Range plus (RR+) method implemented in the OpenStats software package. RR+ is a heuristic method that works on the basis of discretising baseline data into low/normal/high categories. The mutants are then assigned a class based on the reference categories. Finally, Fisher's Exact test applies to specify any significant deviation from the normal category. Here we explain the MP term assignment algorithm for the results from the RR+ framework.

| From the statistical results | From the Annotation Indexer |
|---|--|
| <ol style="list-style-type: none"> Overall effect (do not consider gender) <ul style="list-style-type: none"> if $pvalue_{low} \geq threshold$ & $pvalue_{high} \geq threshold \rightarrow$ Assign no MP term for each $pvalue_{low/high} < threshold$ then assign label 'ABNORMAL', 'INCREASED', 'DECREASED' to the search criteria Remove any Low.INCREASE and High.DECREASE from the labels (II) Apply a similar step to 1 to Male effect (III) Apply a similar step to 1 to Female effect (IV) | Filter for <ol style="list-style-type: none"> Pipeline_stable_id Procedure_group Parameter_stable_id Get available MP terms (I) |
| Find matches between I and II, III, IV <i>but ignore the term Low and High.</i> | |
| Notes <ul style="list-style-type: none"> IF LOW and High MP terms detected then select ABNORMAL term Generally accepted threshold by the IMPC consortium is 0.0001 | |

Categorical data

Categorical data in the IMPC encompasses a range of qualitative measurements such as abnormality in eye, ear, tail and are analysed using Fisher's Exact test implemented in the R package OpenStats. The output MP term for this type of data is a single term Abnormal phenotype if the test is significant. Here we explain the algorithm:

| Step | From the statistical results | From the Annotation Indexer |
|-------------|--|--|
| | <ol style="list-style-type: none"> Overall effect (do not consider gender) <ul style="list-style-type: none"> if $pvalue \geq threshold \rightarrow$ Assign no MP term if $pvalue < threshold$ search for the MP term (II) Apply a similar step to 1 to Male effect (III) Apply a similar step to 1 to Female effect (IV) | Filter for <ol style="list-style-type: none"> Pipeline_stable_id Procedure_group Parameter_stable_id Get available MP terms (I) |
| | Find matches between I and II, III, IV | |
| Note | Generally accepted threshold by the IMPC consortium is 0.0001 | |

Schematic view of the IMPC-AP

