

Rümmler, Arne  
Chair of Geoinformatics, TU Dresden

# How to evaluate and reduce the complexity of geospatial provenance graphs?

01.09.2020

# PROVENANCE

- Provenance describes the origination history of data

# PROVENANCE

- Provenance describes the origination history of data
- Workflow graphs are a widely known form of provenance

# PROVENANCE

- Provenance describes the origination history of data
- Workflow graphs are a widely known form of provenance
- Provenance data can be standardized (e.g. PROV-O)

# PROVENANCE

- Provenance describes the origination history of data
- Workflow graphs are a widely known form of provenance
- Provenance data can be standardized (e.g. PROV-O)
- Standardized provenance graphs can be linked

# PROVENANCE

- Provenance describes the origination history of data
- Workflow graphs are a widely known form of provenance
- Provenance data can be standardized (e.g. PROV-O)
- Standardized provenance graphs can be linked
- Provenance graphs can become indefinitely large

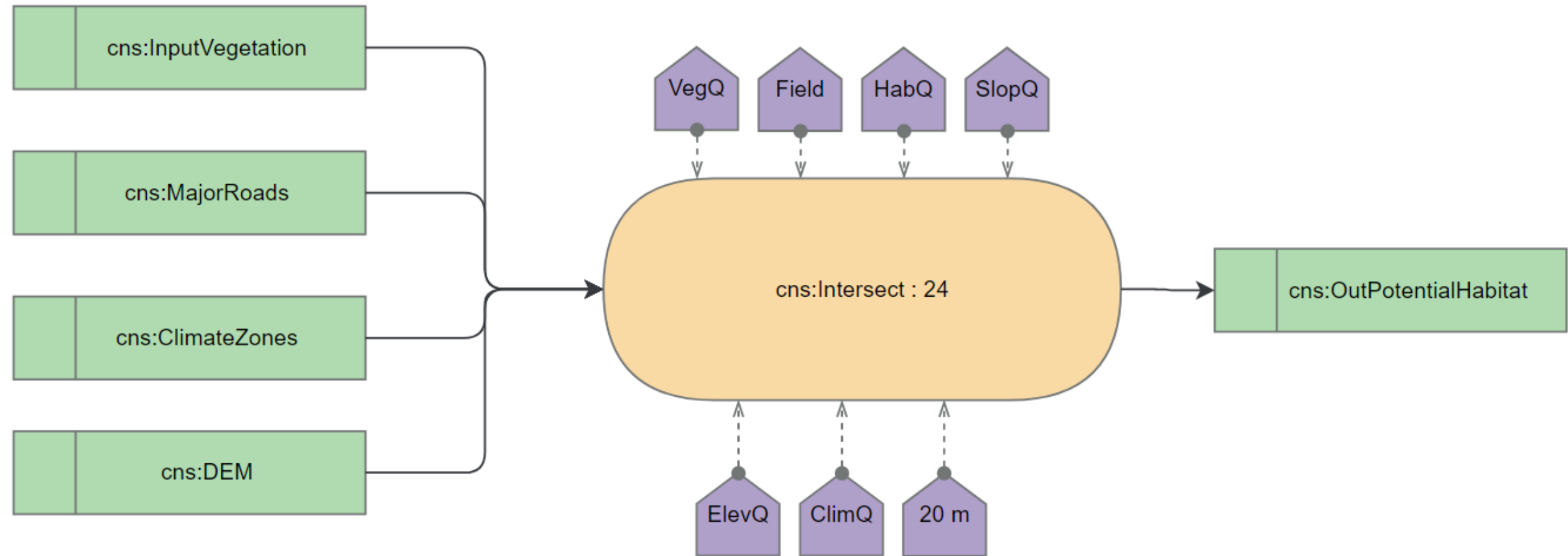
# PROVENANCE

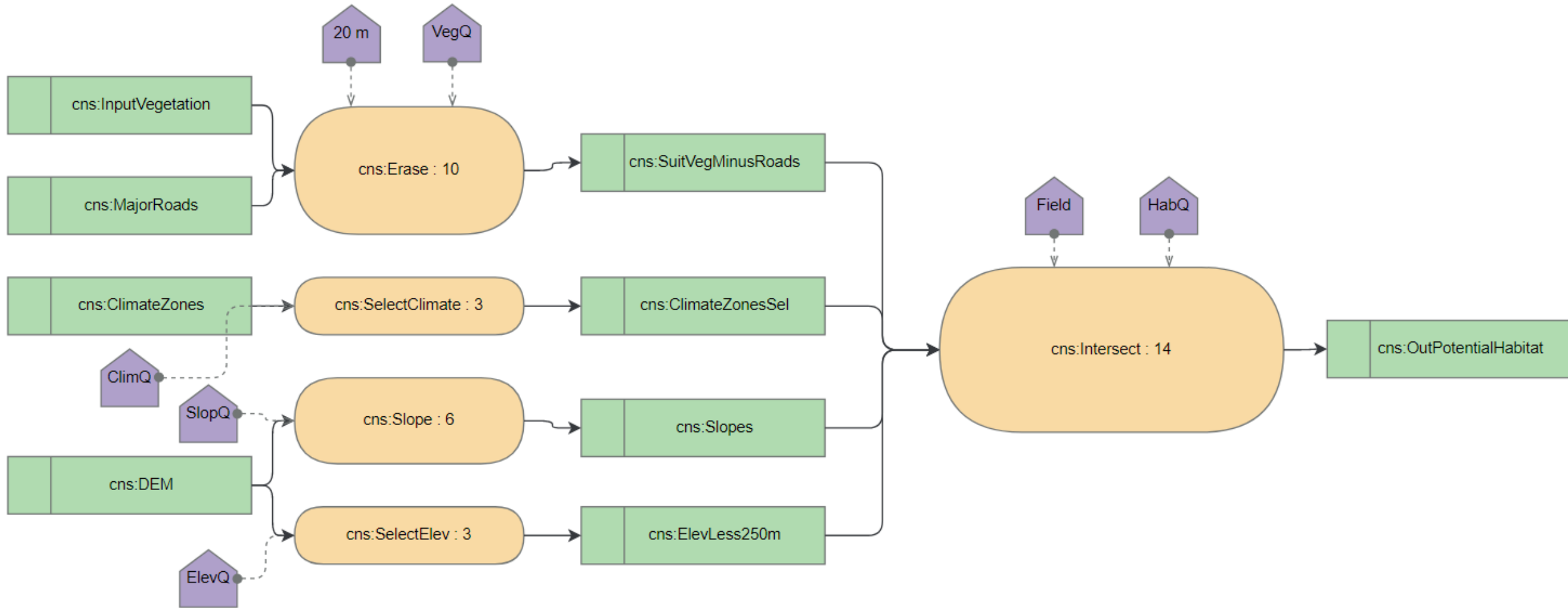
- Provenance describes the origination history of data
- Workflow graphs are a widely known form of provenance
- Provenance data can be standardized (e.g. PROV-O)
- Standardized provenance graphs can be linked
- Provenance graphs can become indefinitely large

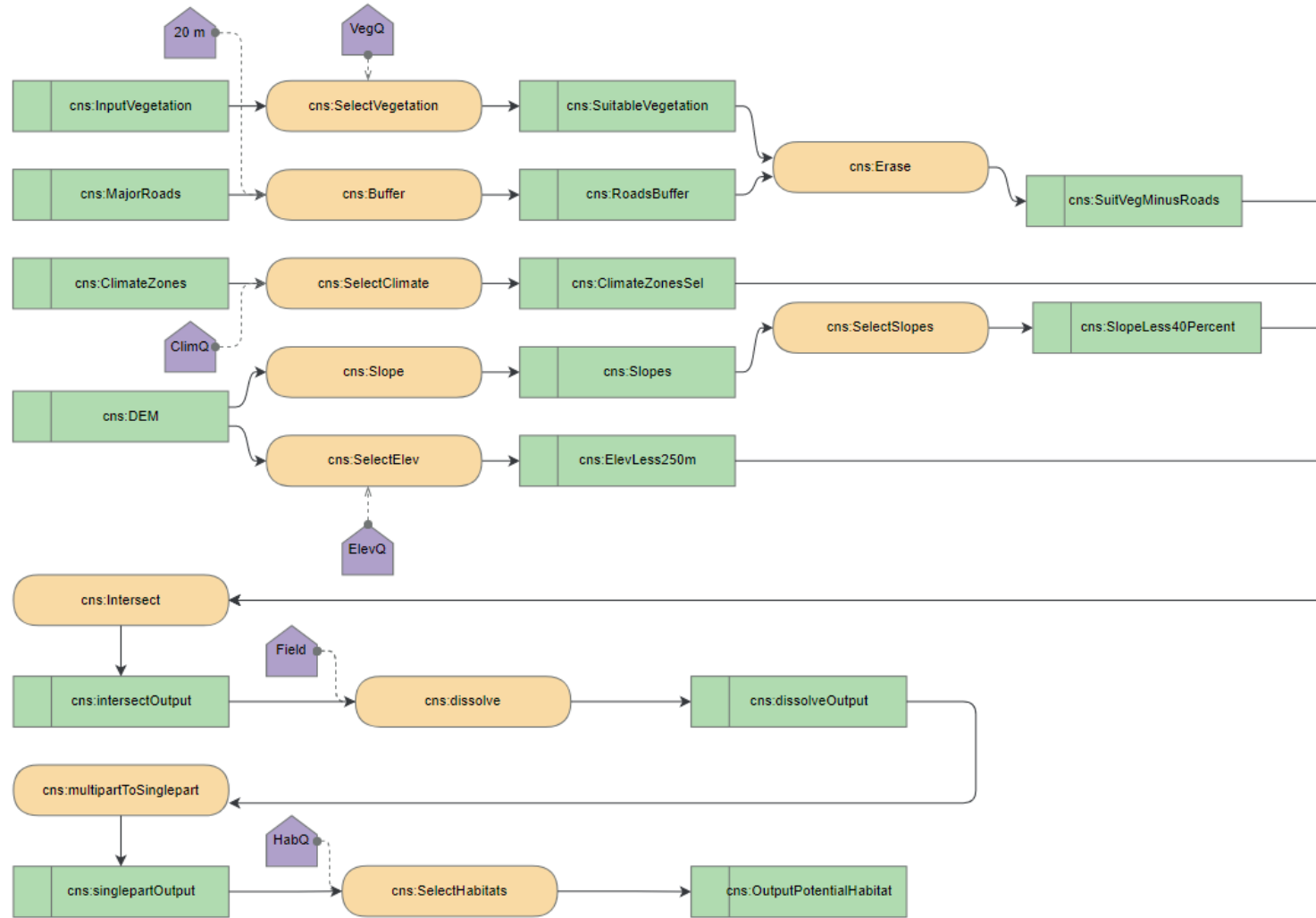
**How to make use of such provenance data?**

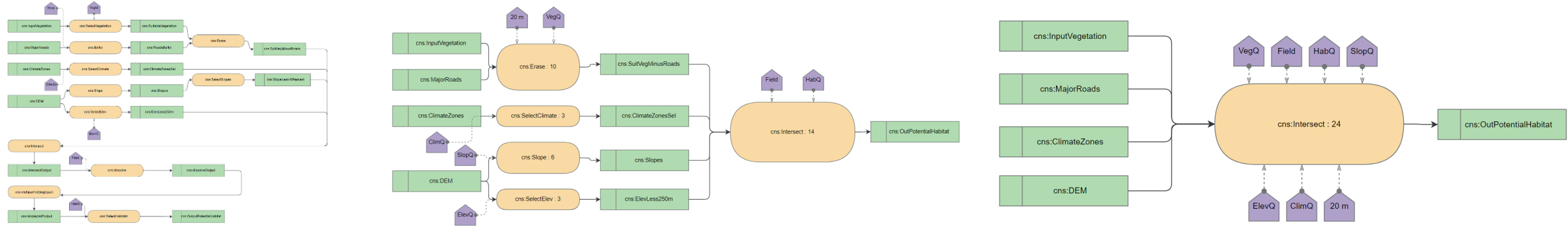
Leverage standardized provenance data to enhance the understandability of the origination history of geospatial data.





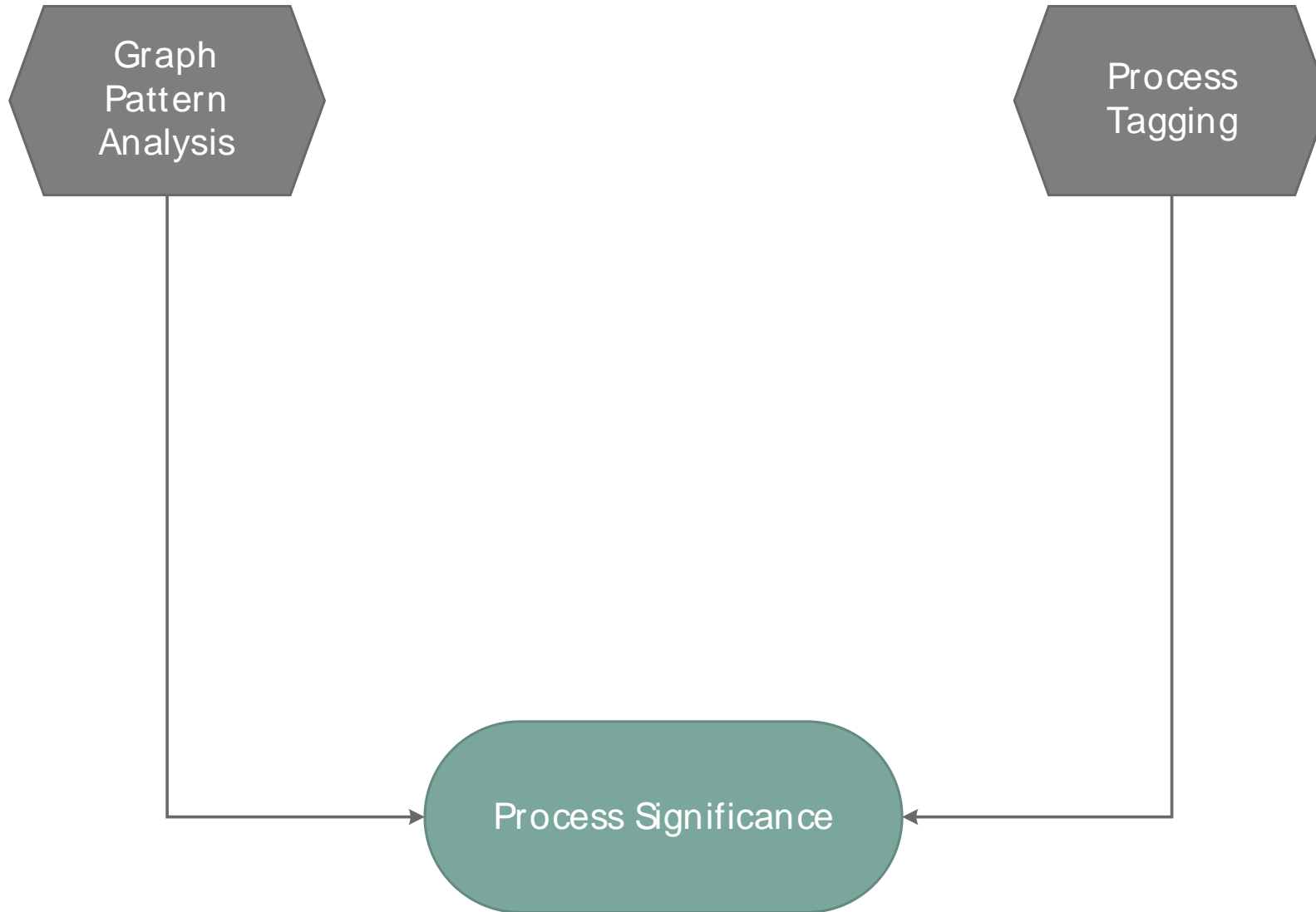


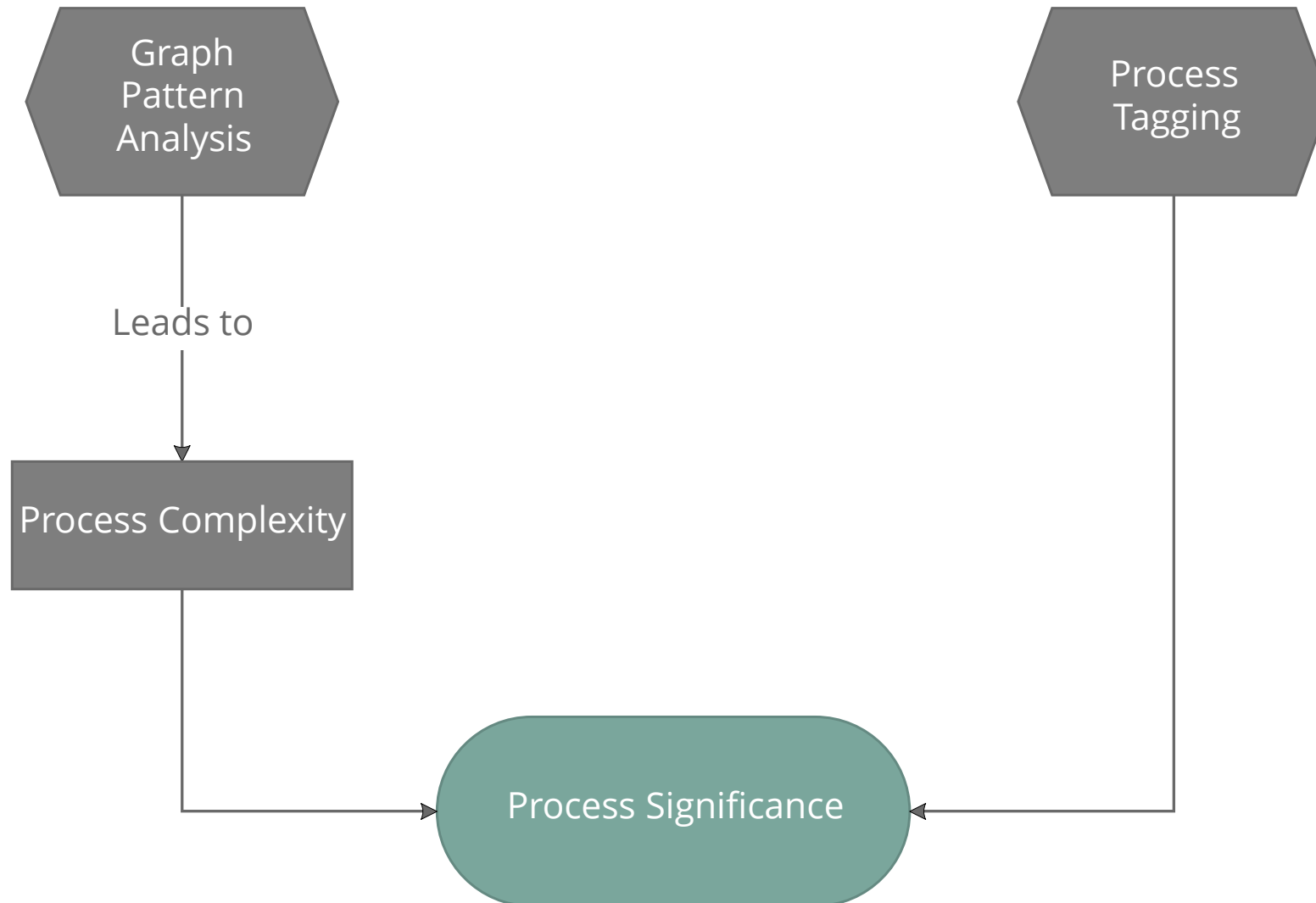


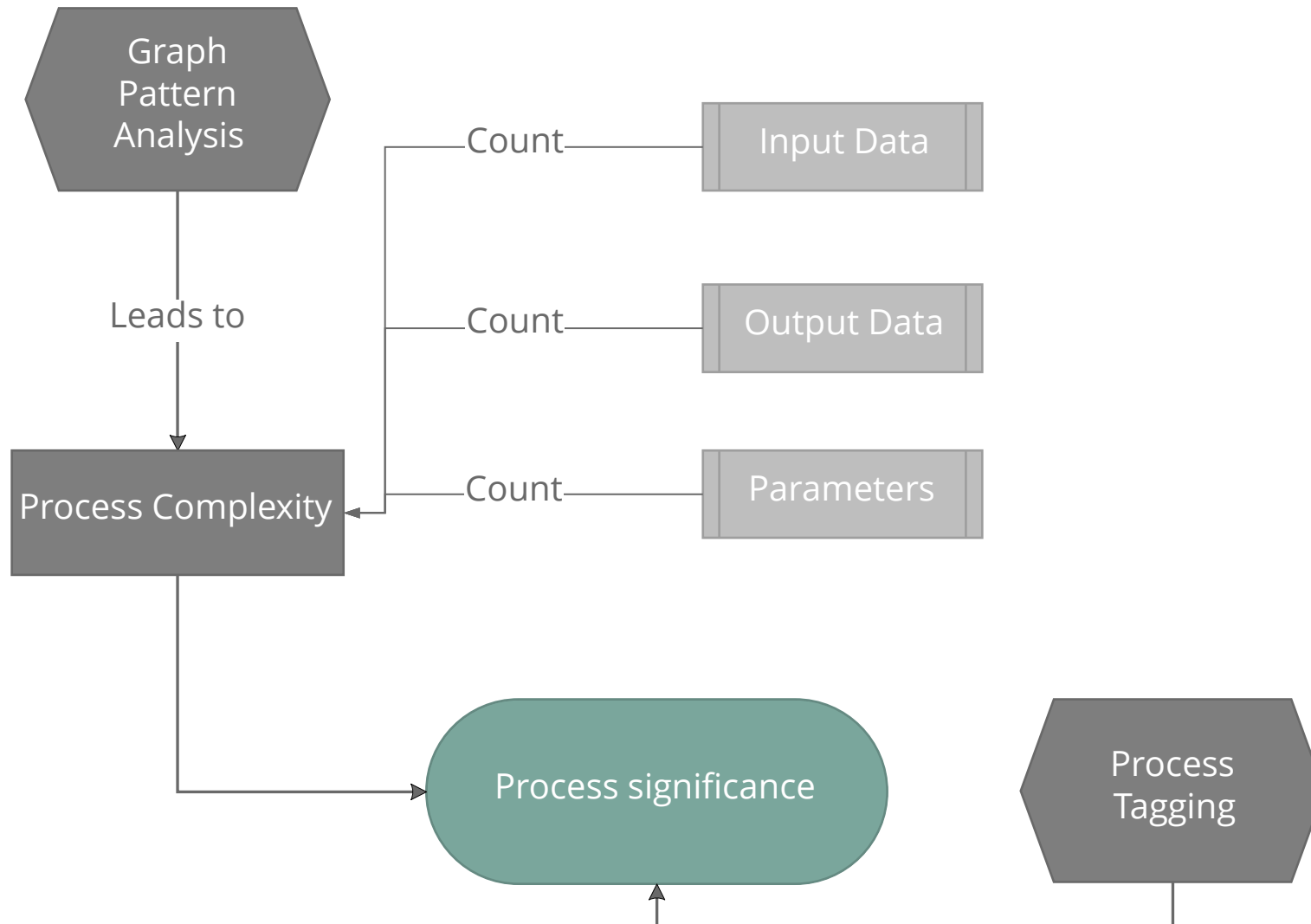


How to get from the left to the right?

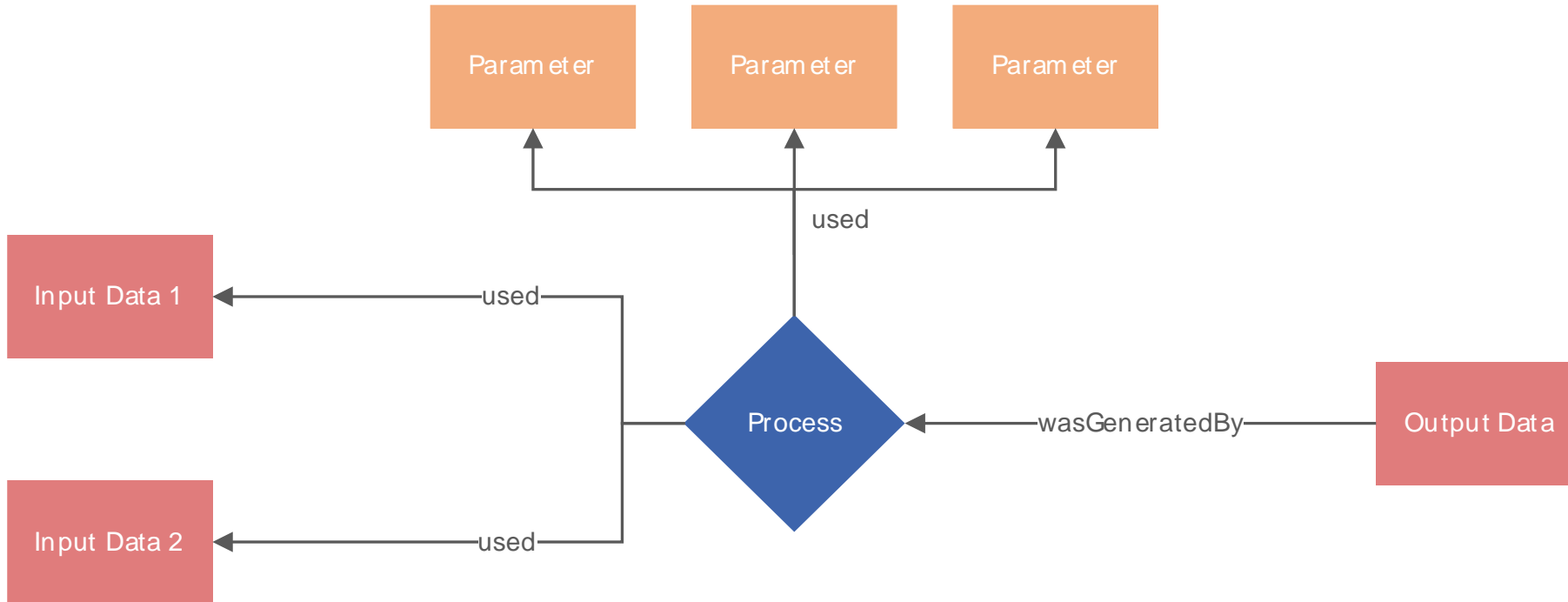
-> **Process-significance based generalization**





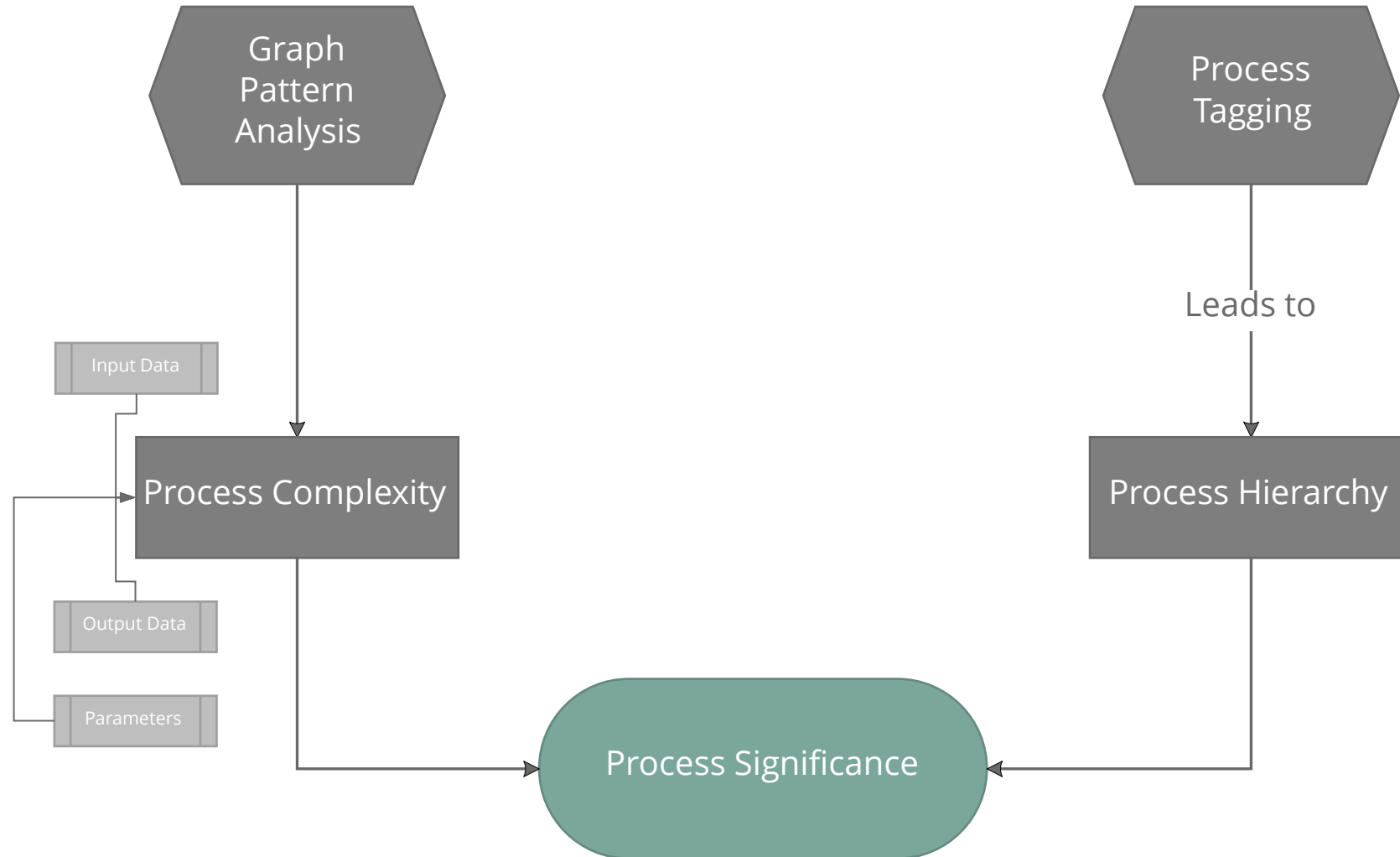


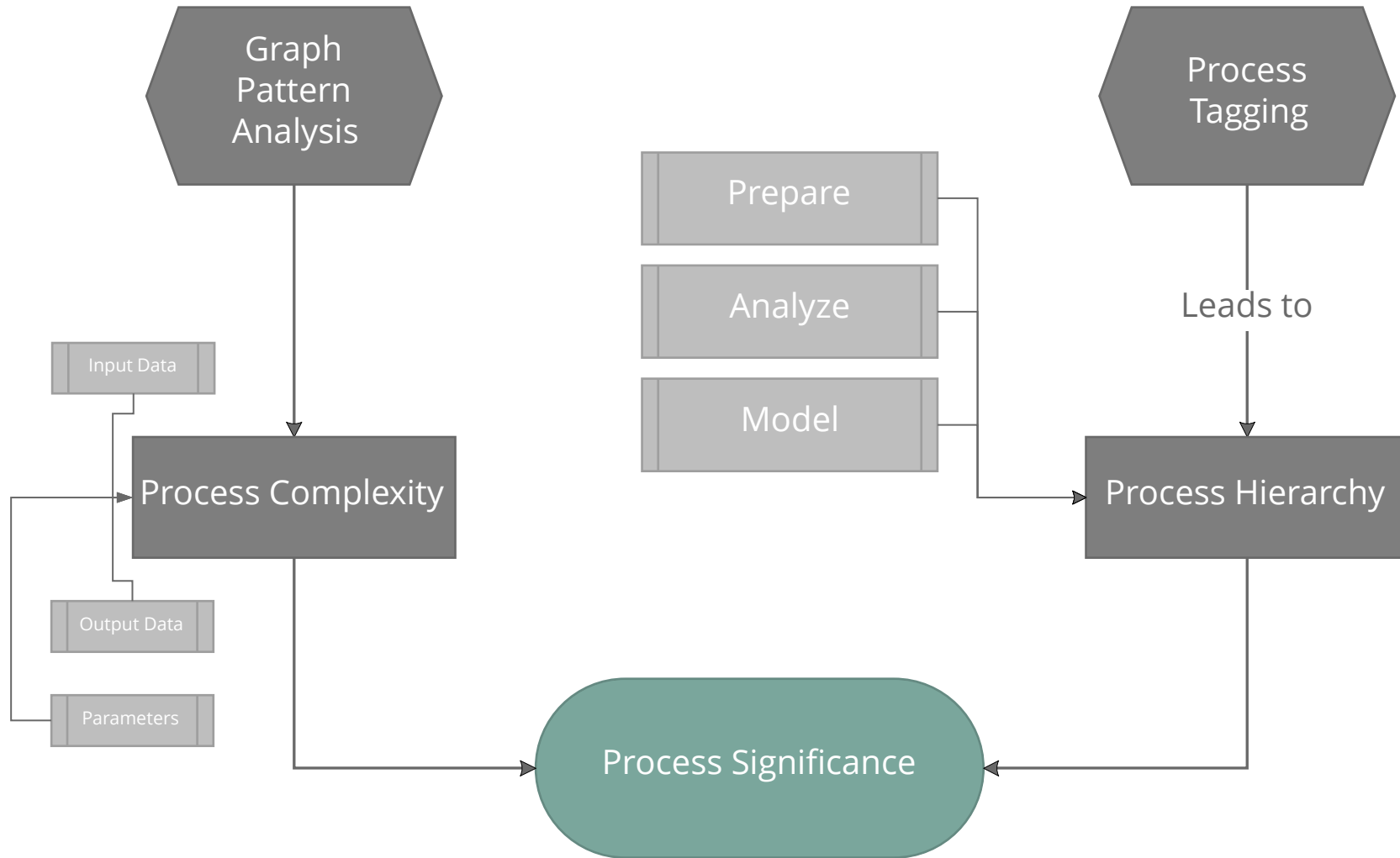
## Process complexity analysis by graph pattern



- 2 Inputs, 3 Parameter, 1 Output -> Complexity  $C = 2 + 3 + 1 = 6$
- Alternative complexity calculation possible, e.g.  $\sqrt{\text{params}}$
- Further patterns can be included

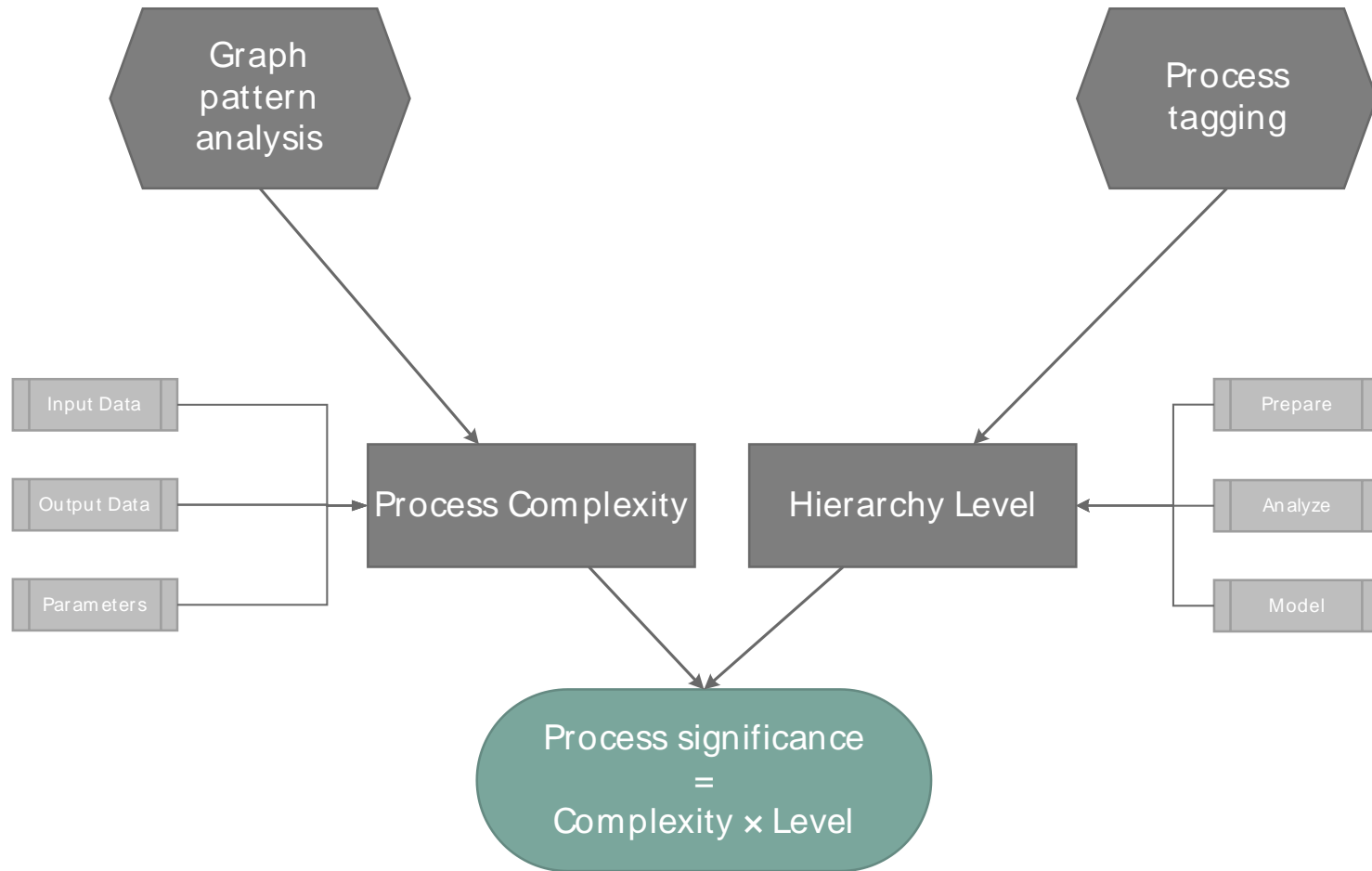






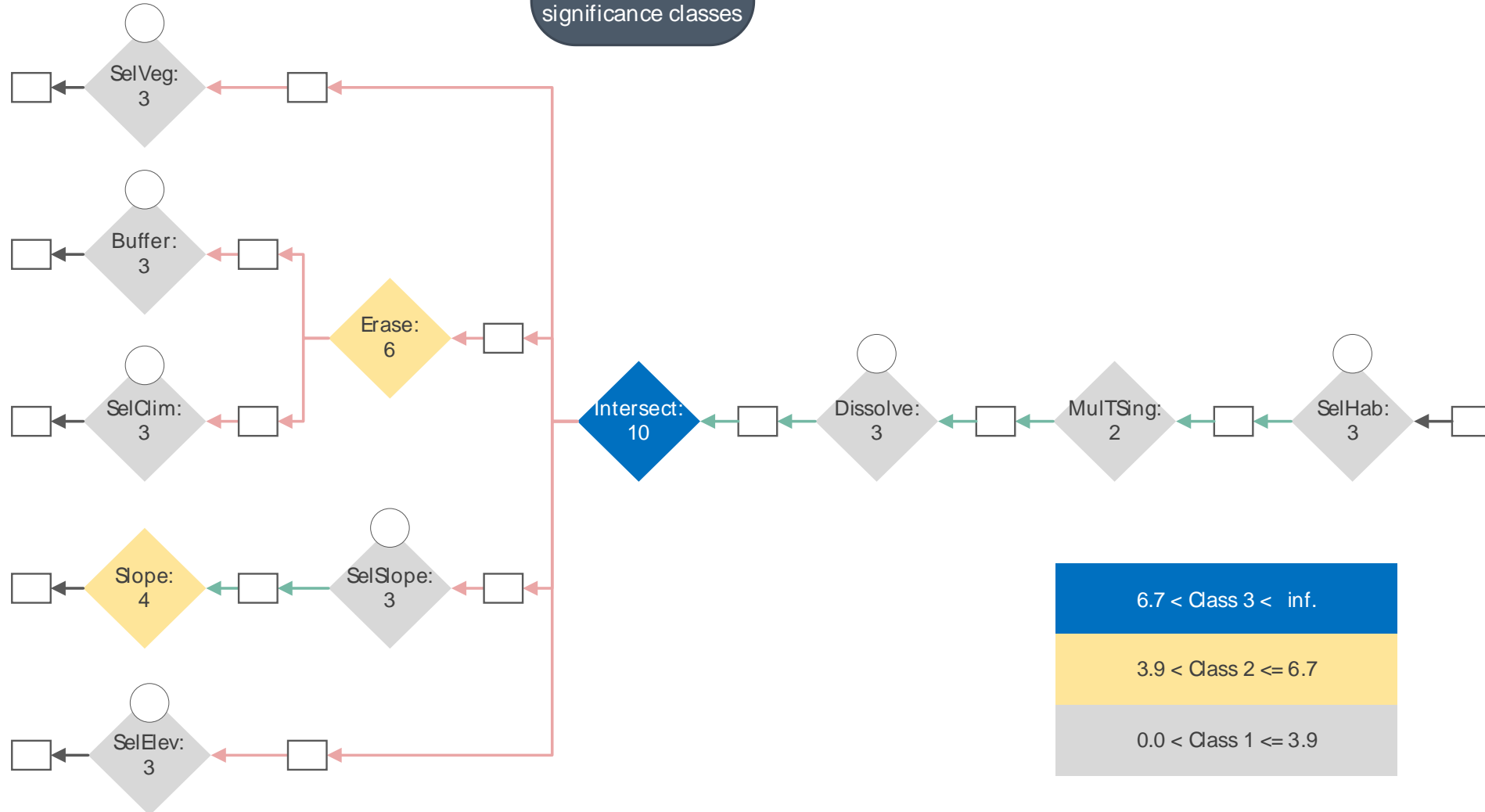
Prepare	Analyze	Model
<ul style="list-style-type: none"> <li>▪ Selection</li> <li>▪ Transformation</li> <li>▪ Reclassification</li> <li>▪ Geometric Computation</li> <li>▪ Visualization</li> </ul>	<ul style="list-style-type: none"> <li>▪ Neighborhood Analysis</li> <li>▪ Map Overlay</li> <li>▪ Network</li> <li>▪ Statistics</li> </ul>	<ul style="list-style-type: none"> <li>▪ Model (back-box-modules, neural networks, ... )</li> </ul>
1	2	3

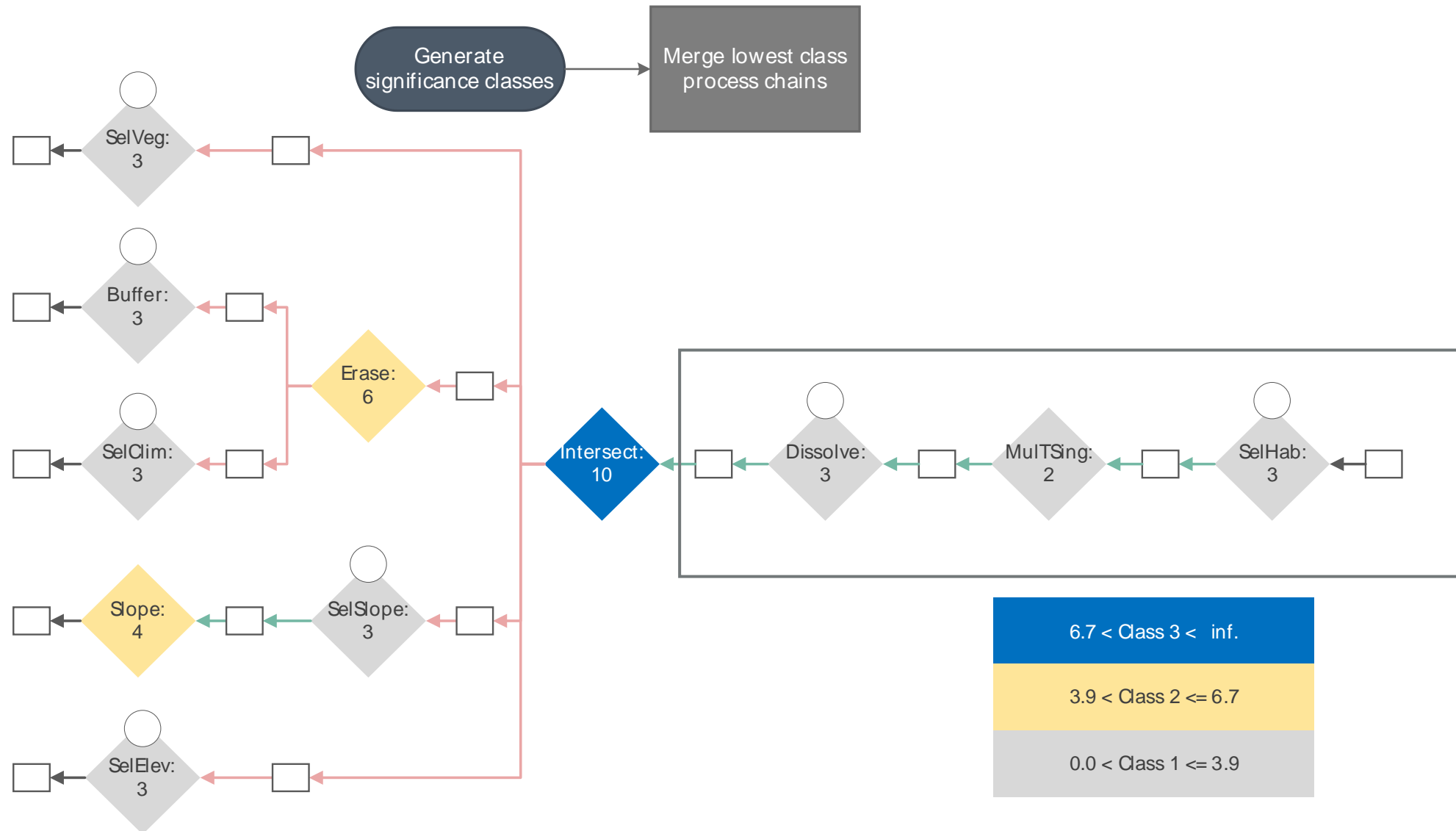
- Subcategories derived from Albrecht (Albrecht, J. Universal analytical GIS operations: a task-oriented systematization of data structure-independent GIS functionality 1998.)
- Value serves a multiplier

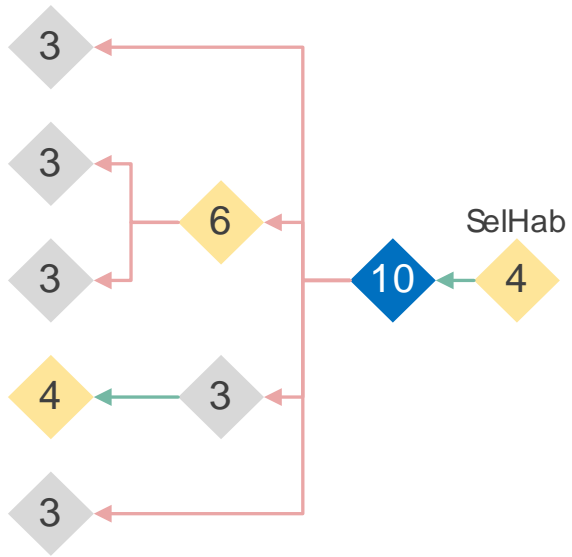


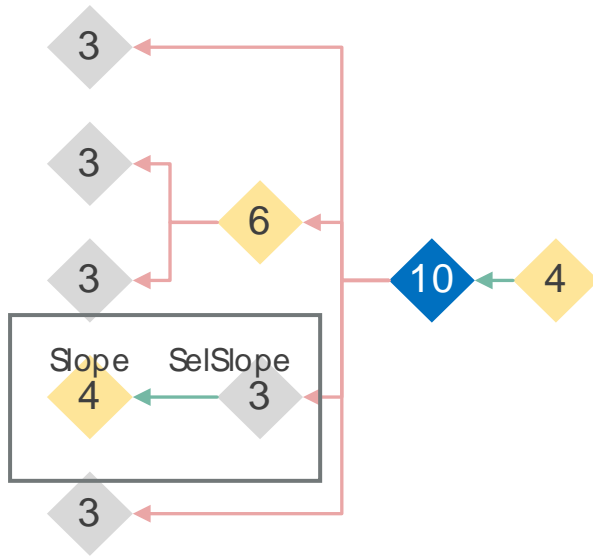
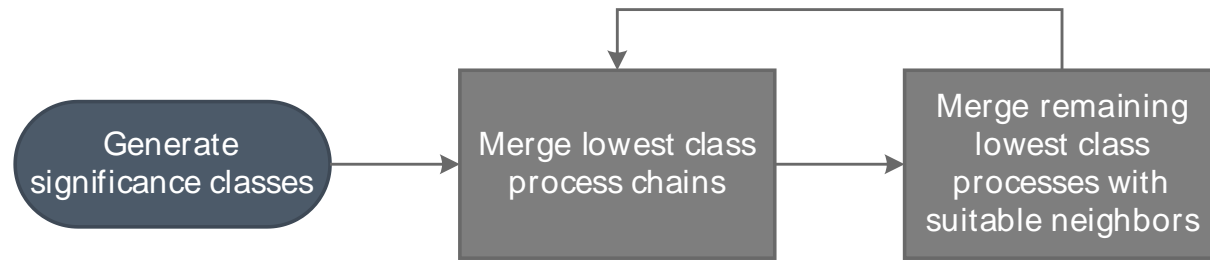
# Graph Generalization

Generate  
significance classes

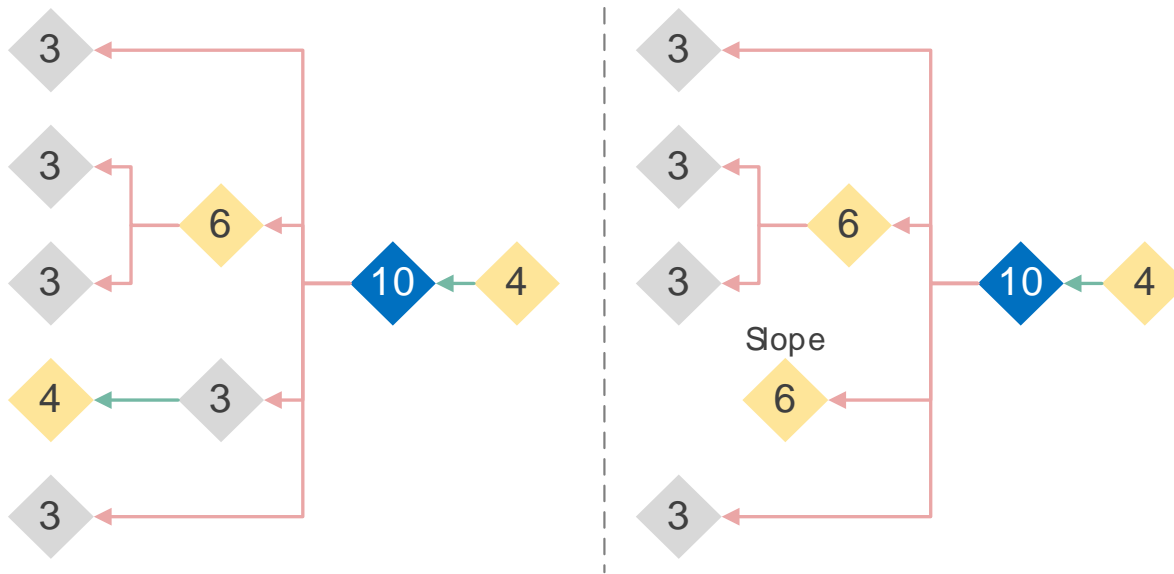
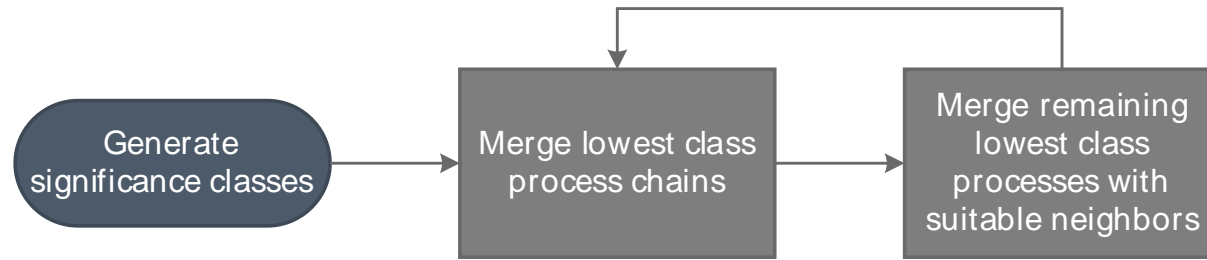


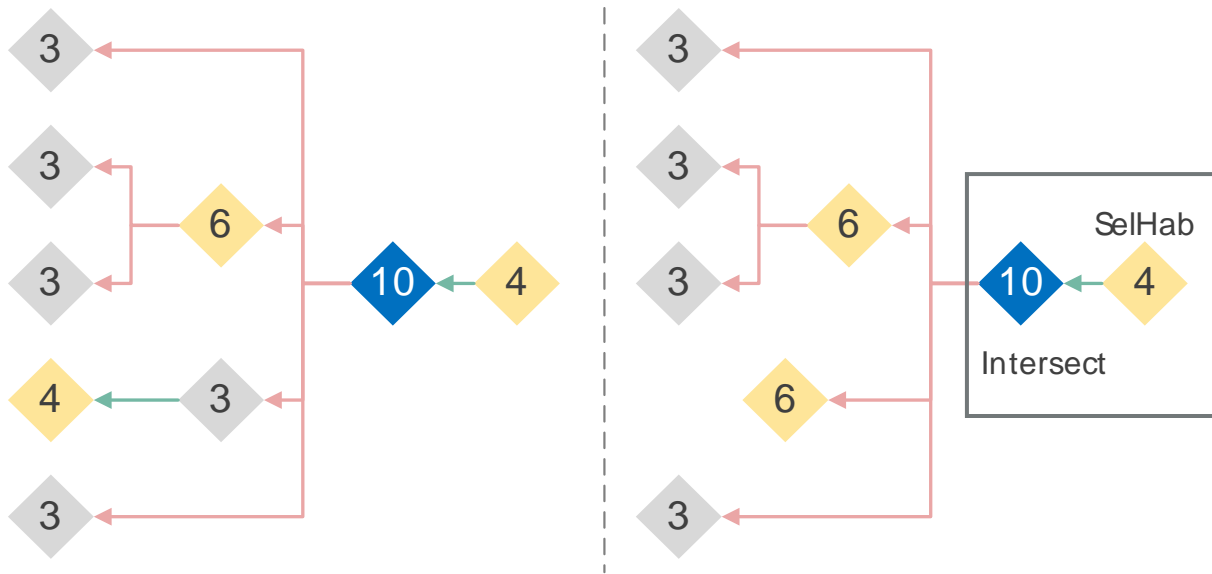
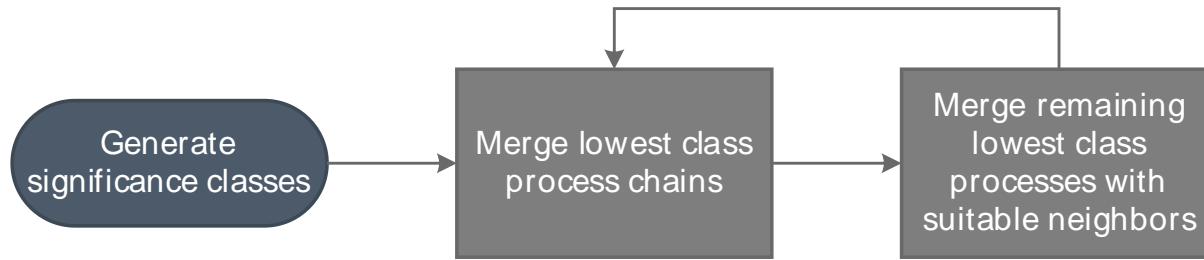


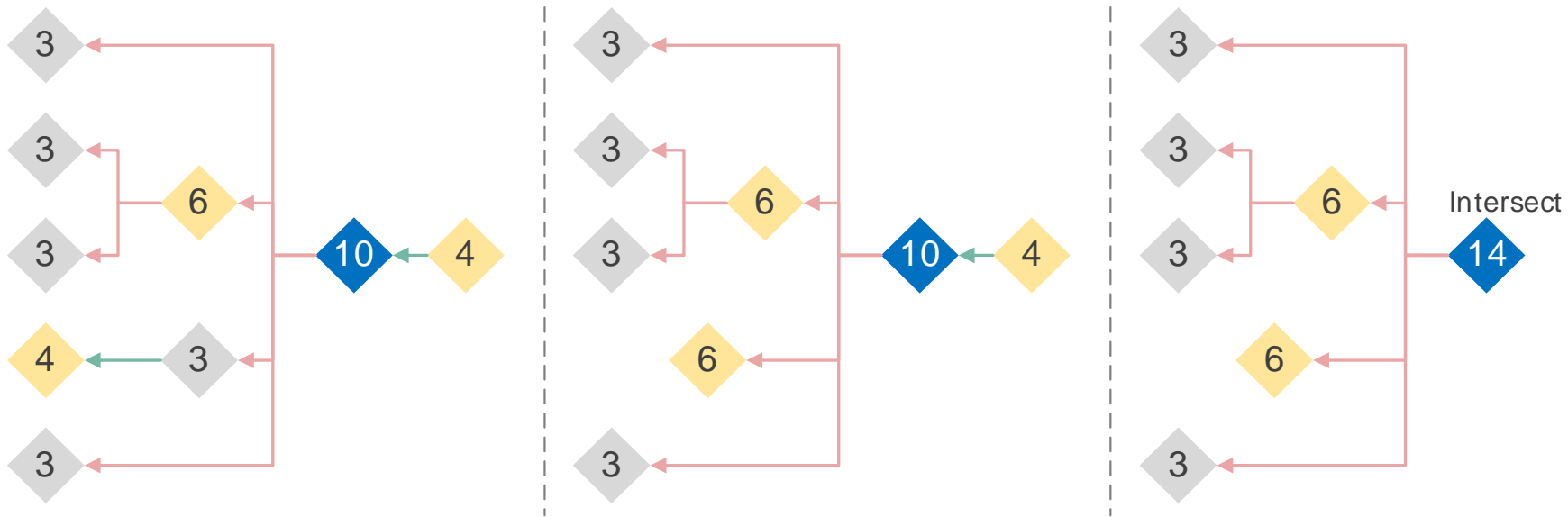
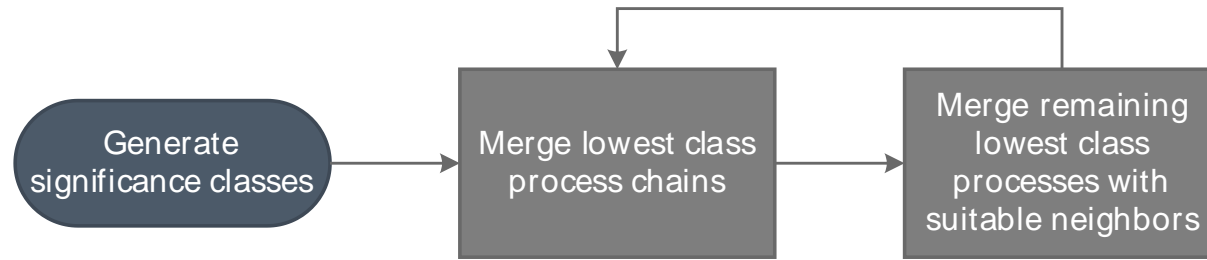


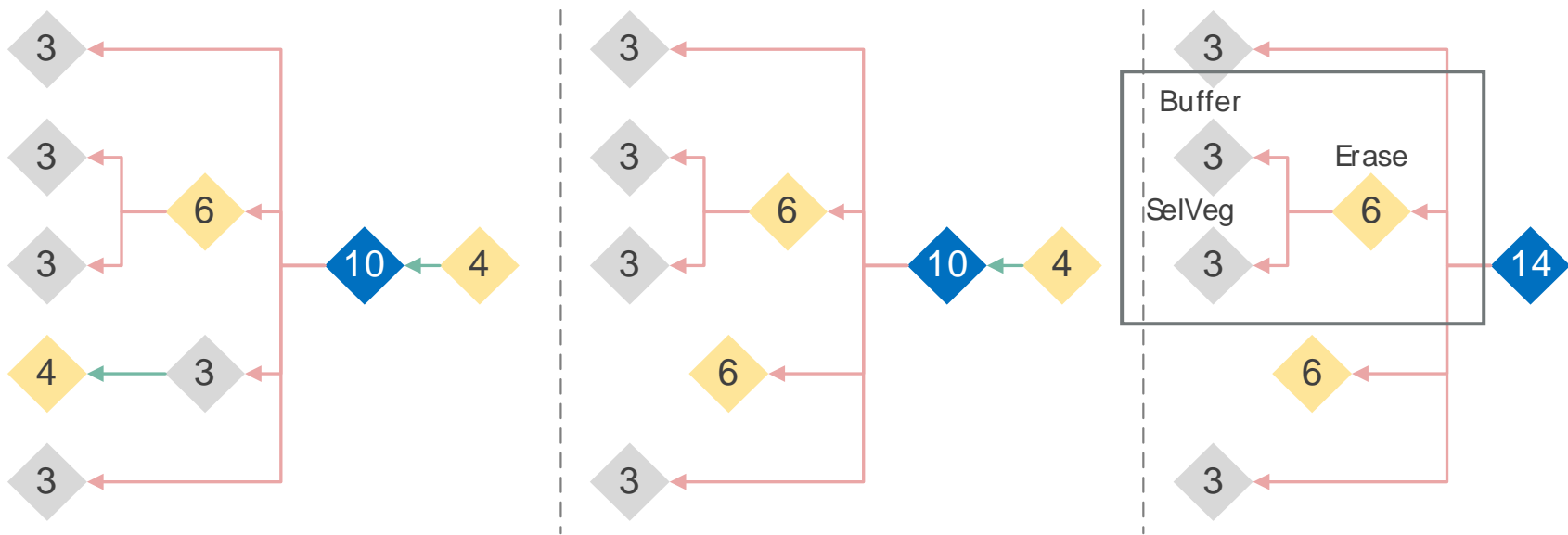
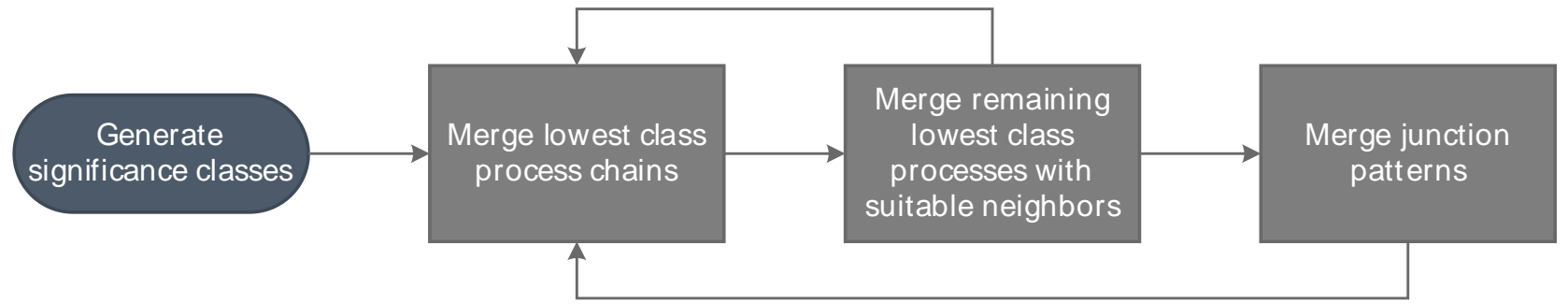


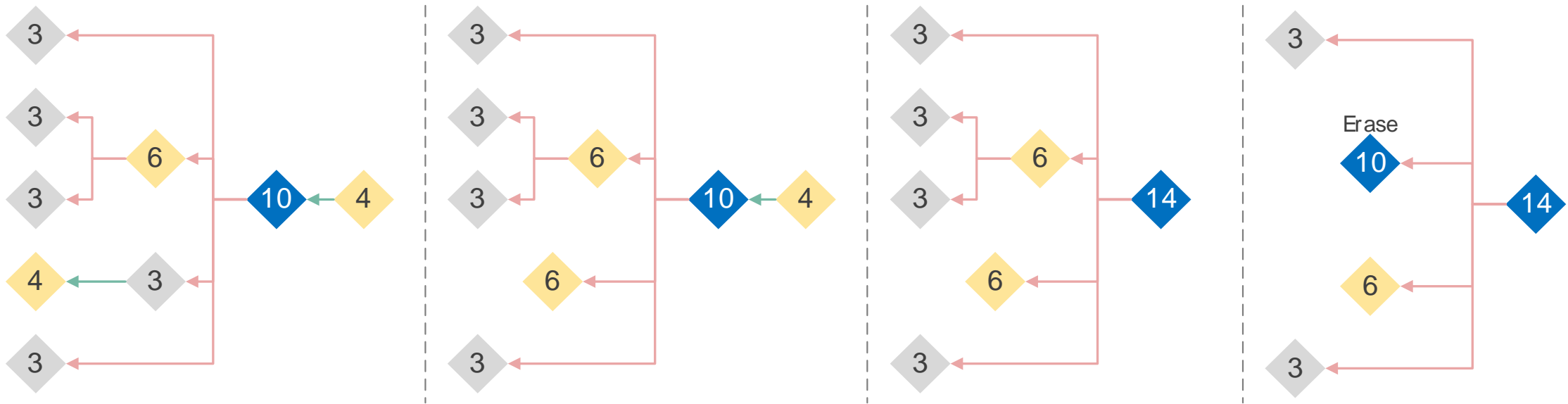
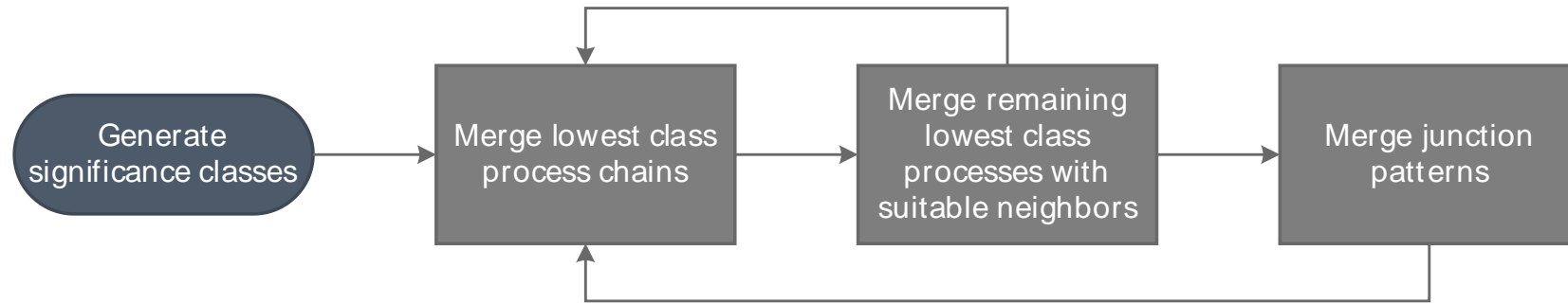


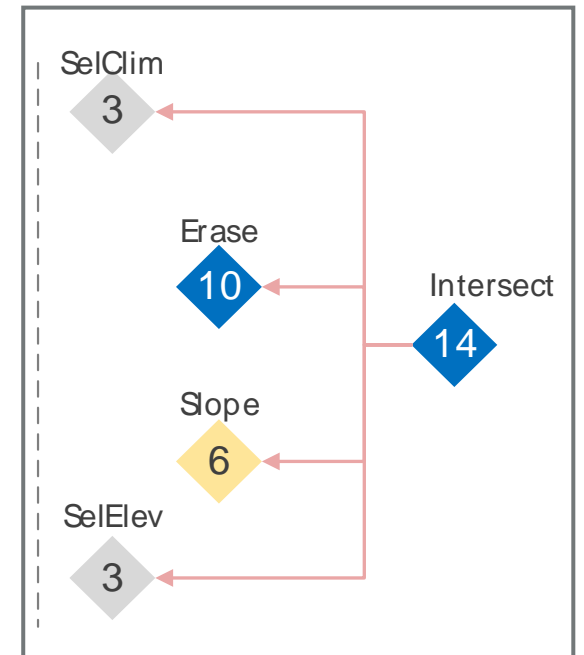
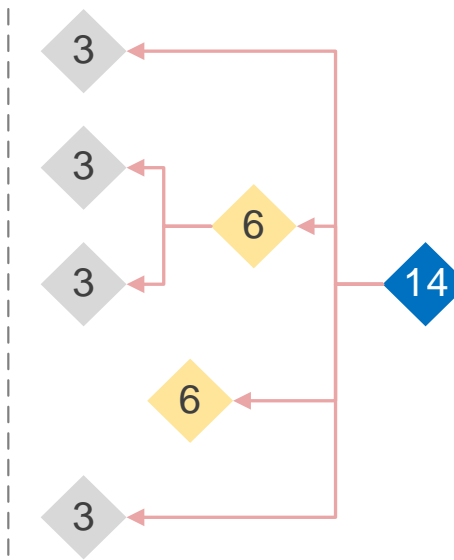
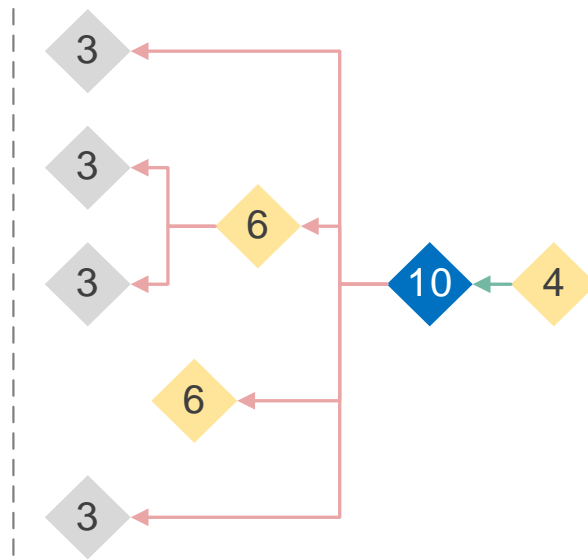
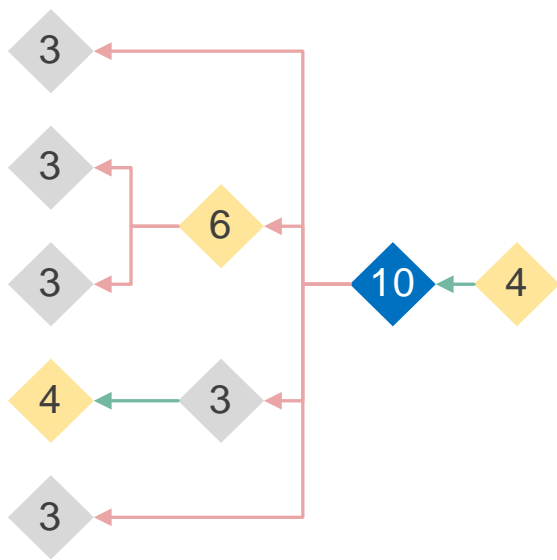
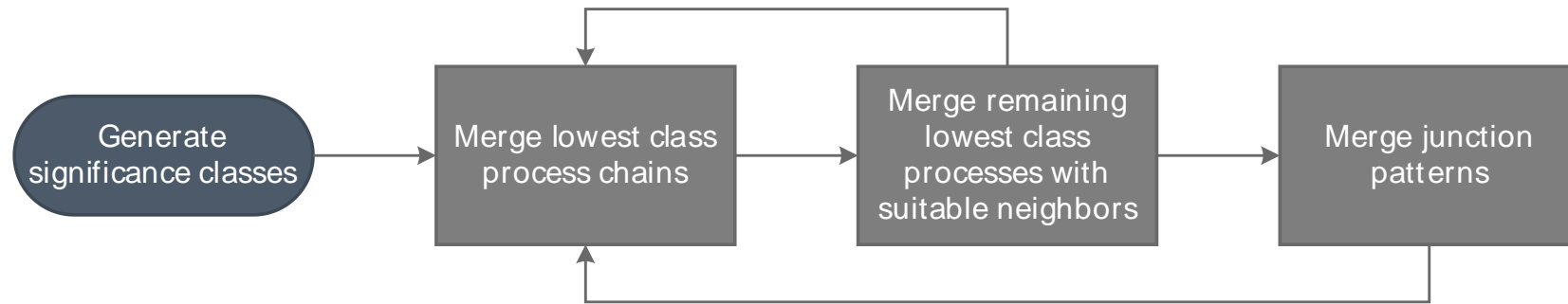


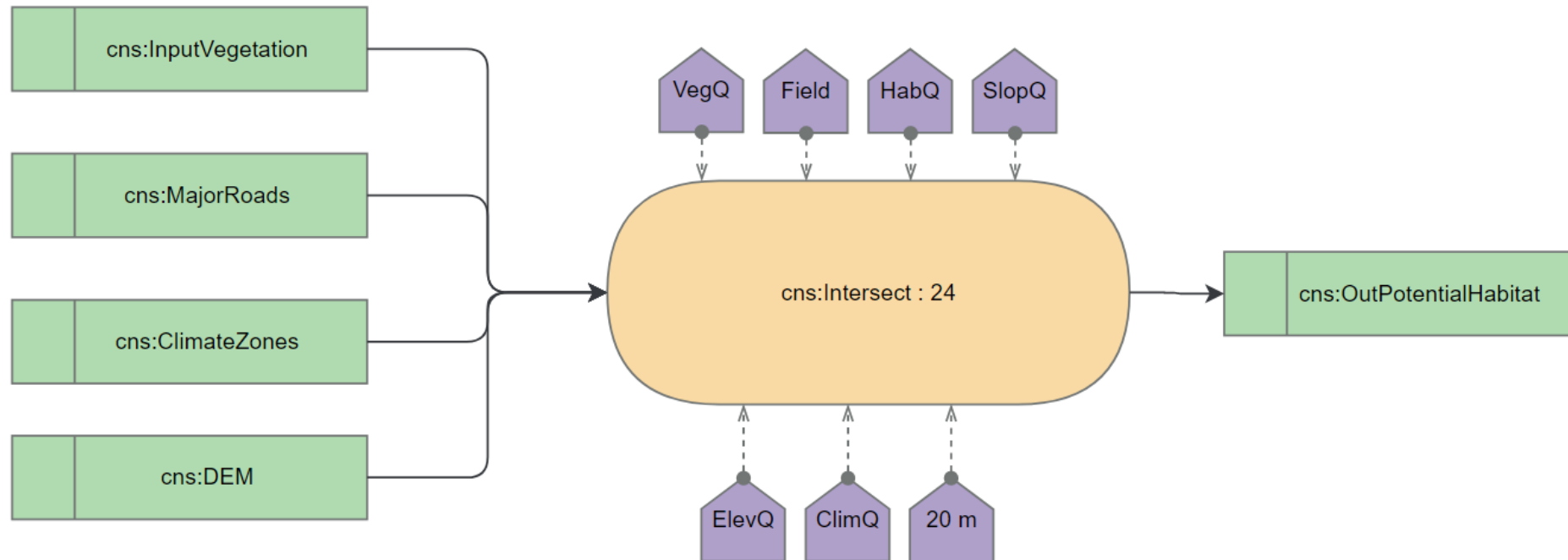
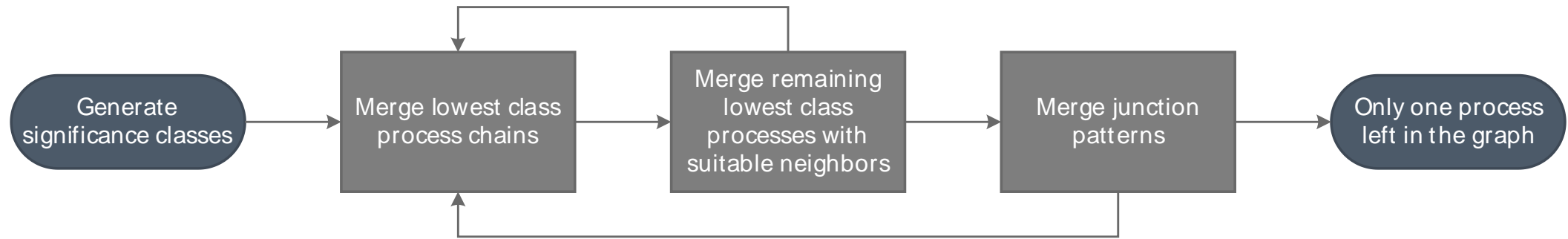






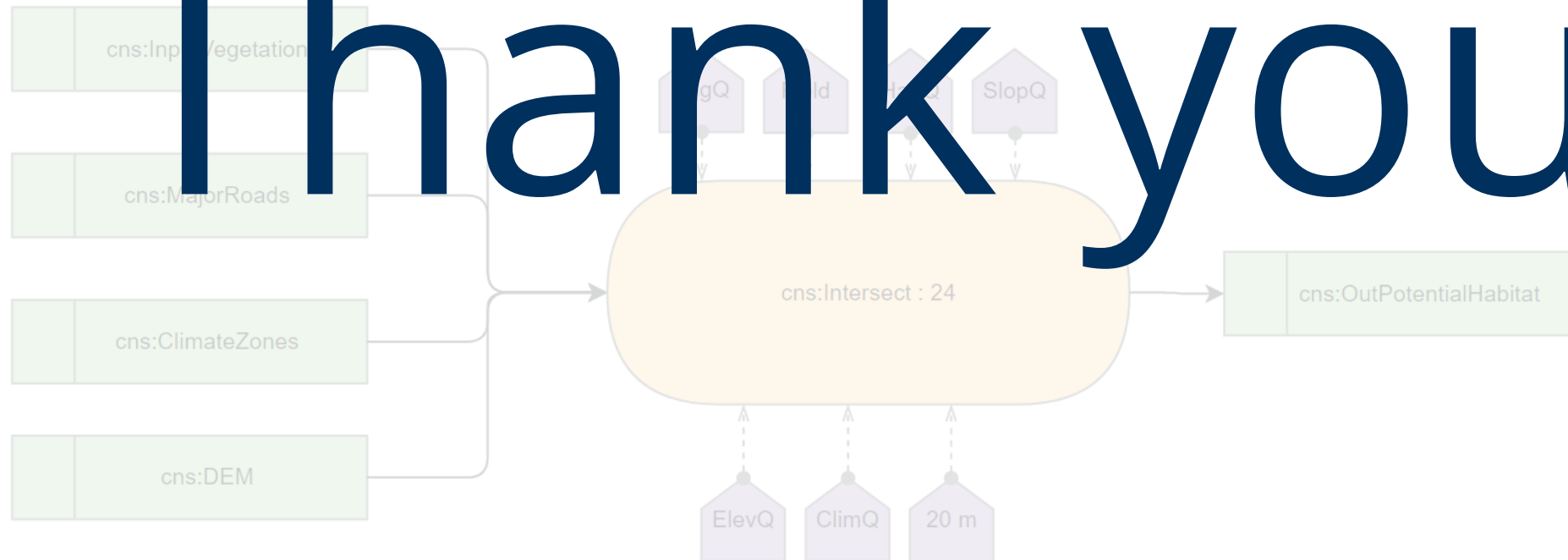




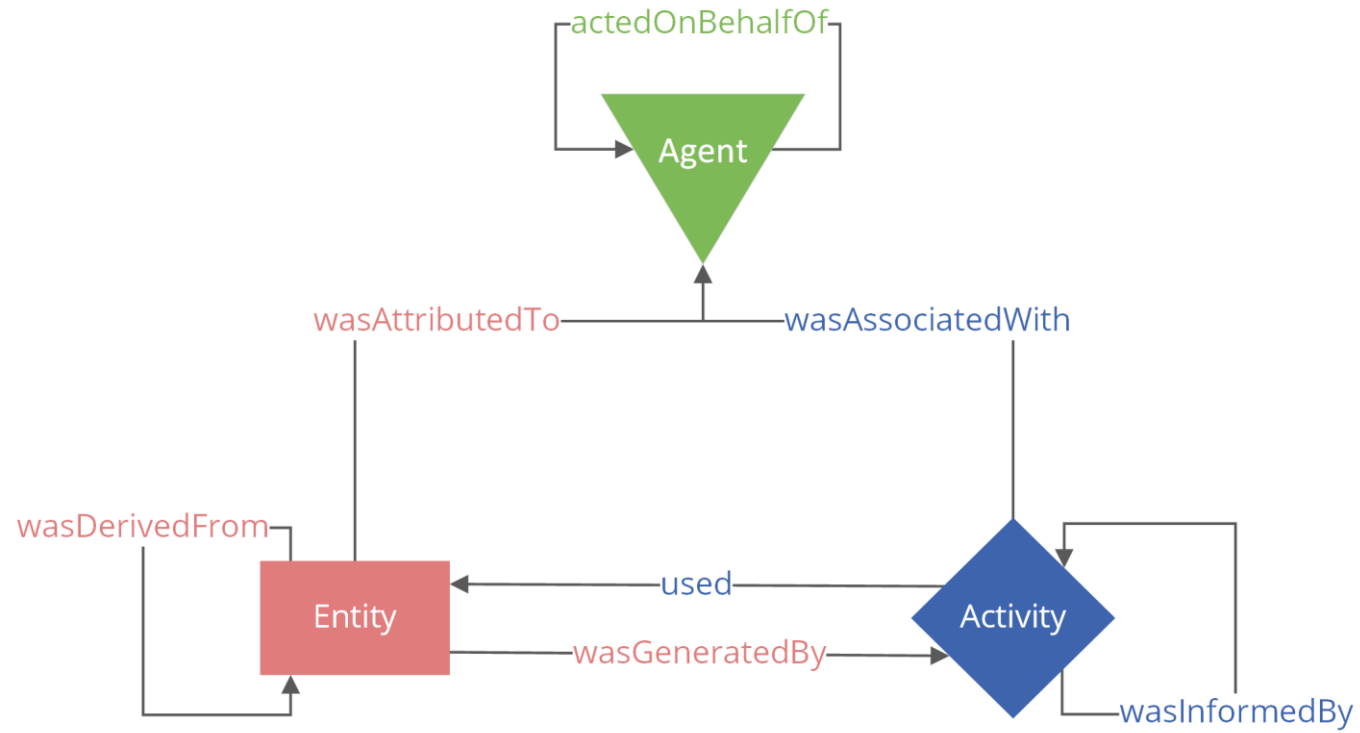


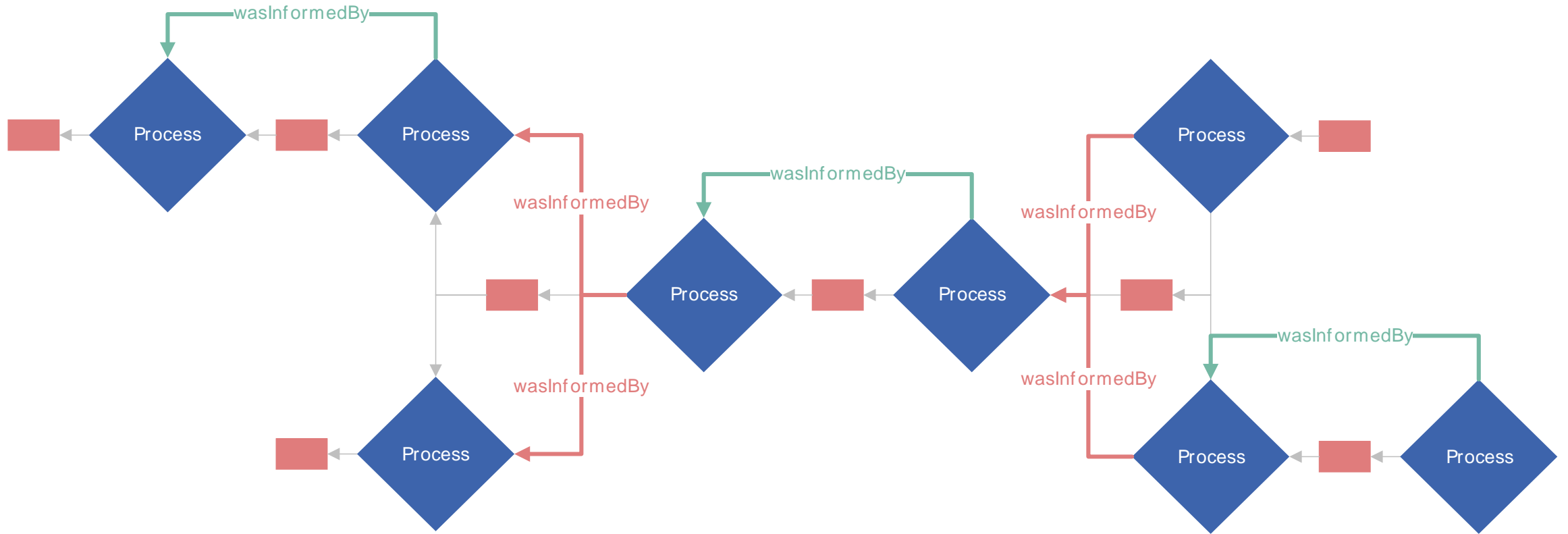


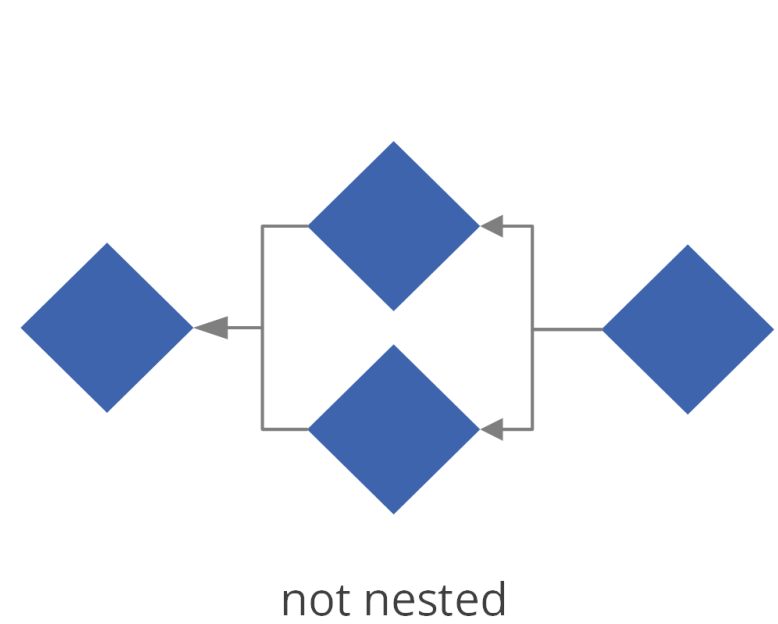
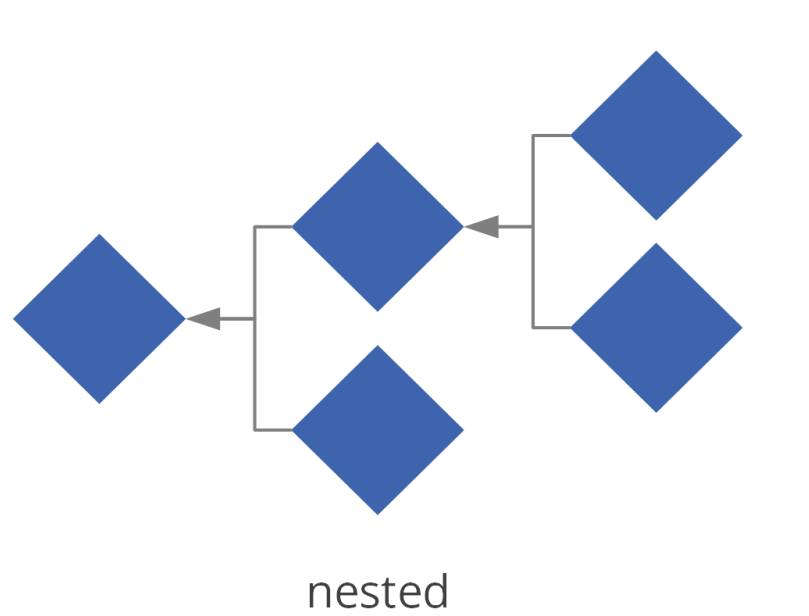
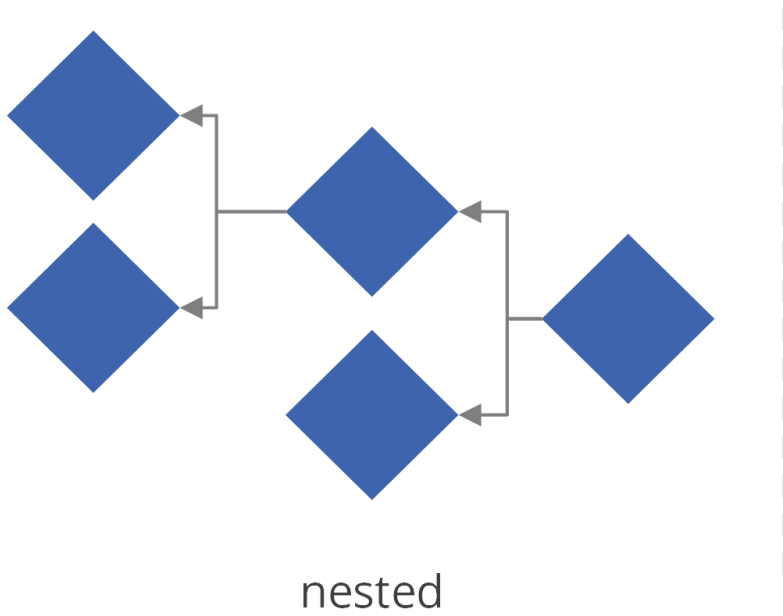
# Thank you





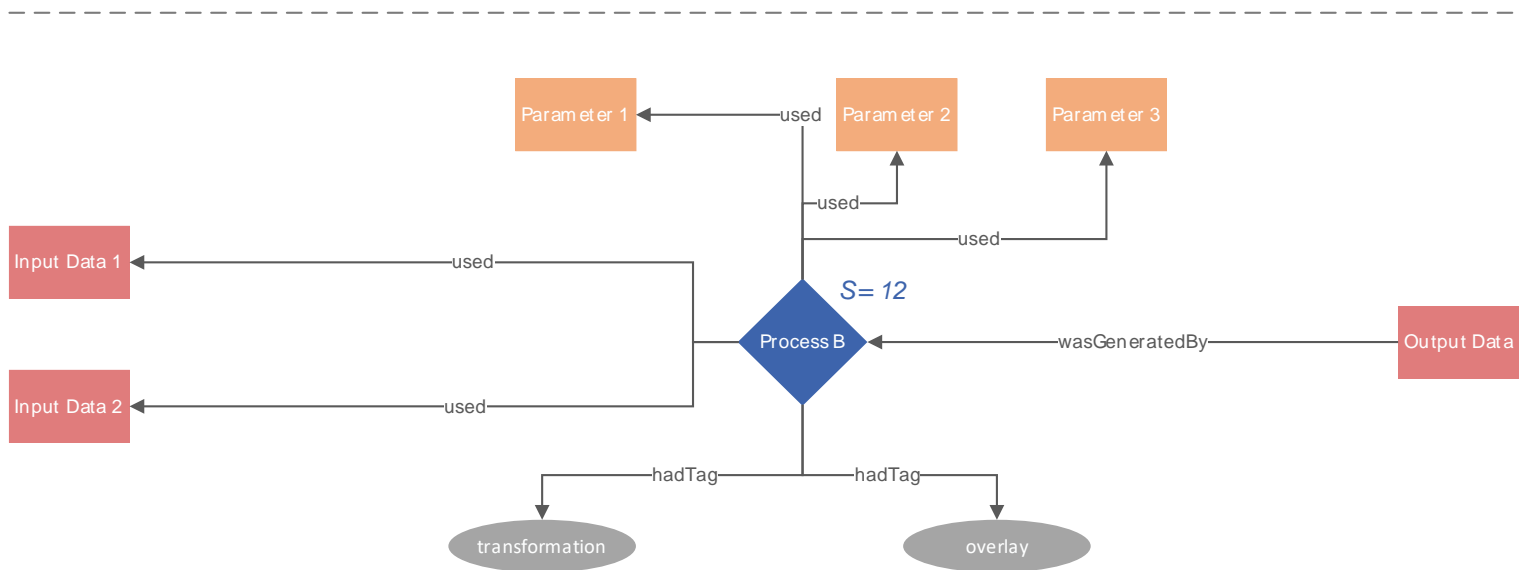
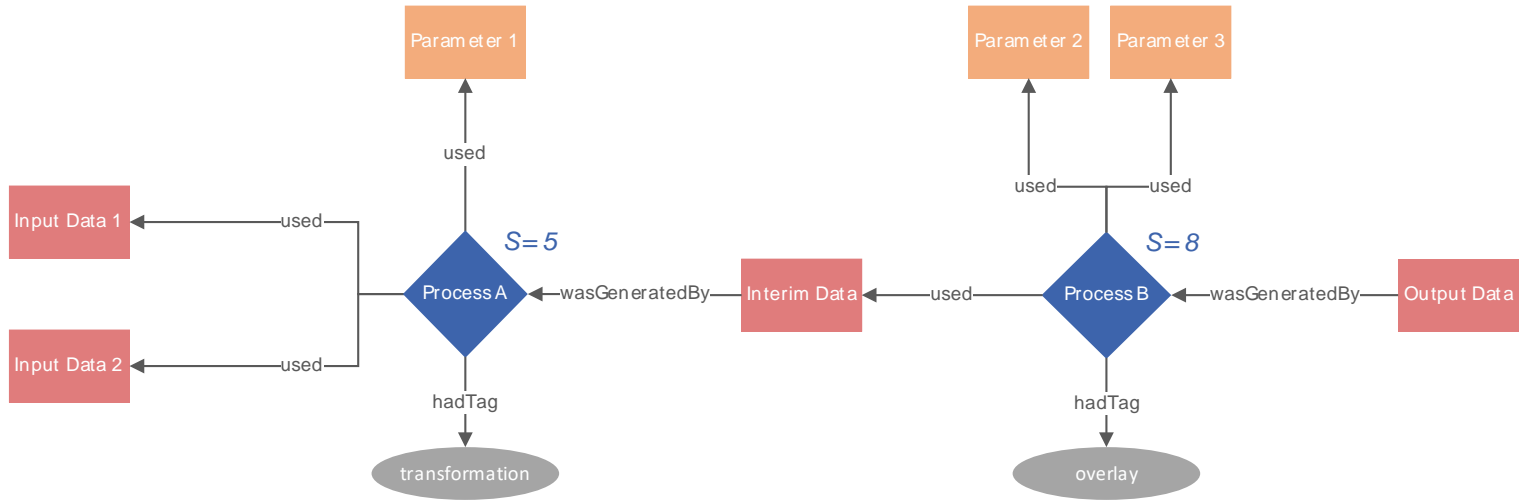


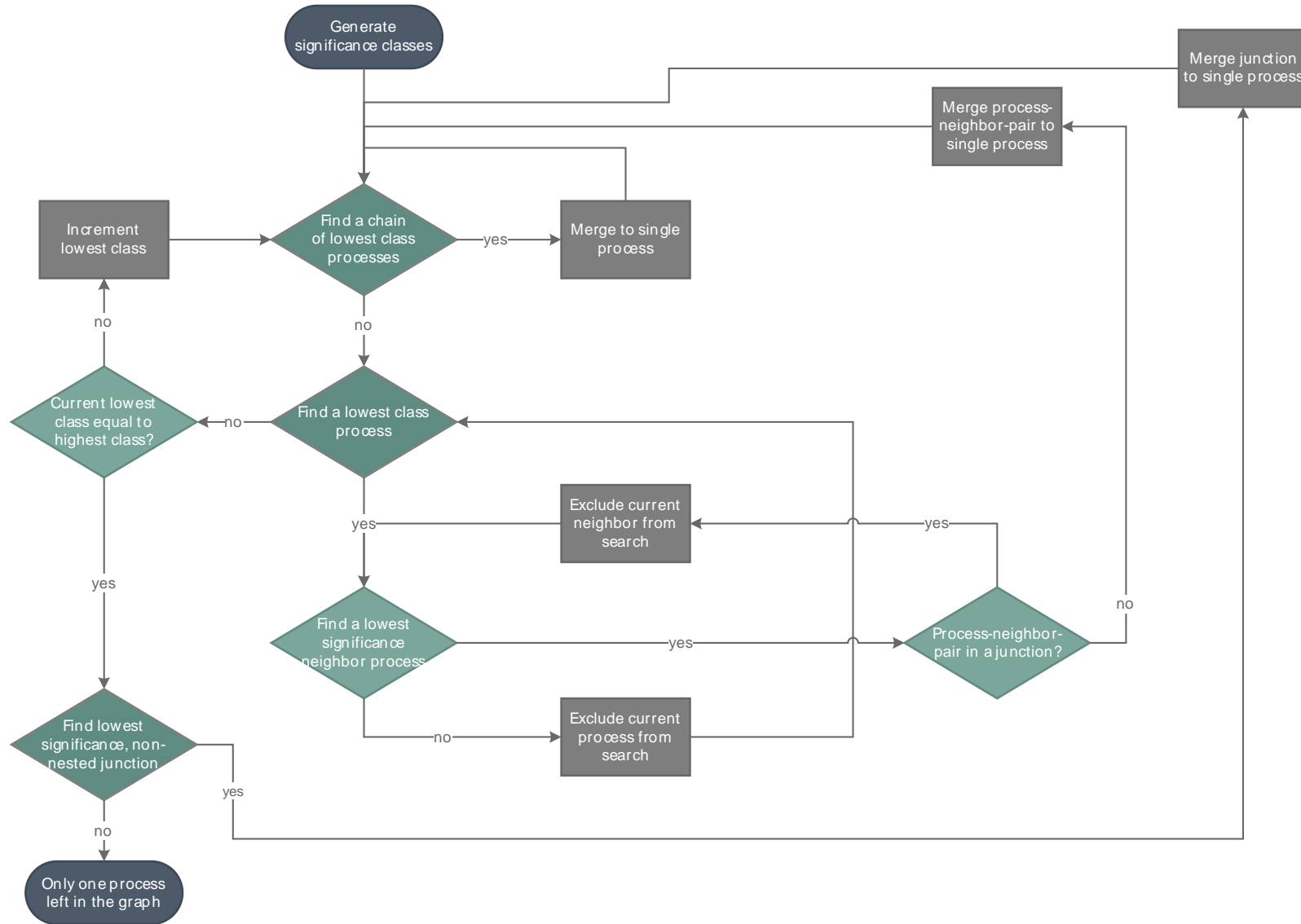




## Merging Rules

- Any interim data between the merged Processes is ignored. Interim data is data that is at the same time input data and output data (only of the Processes that are getting merged).
- The input data of the new Process equals the source data of the merged units. Source data is input data that is no interim data
- The output data of the new Process is equal to the final data of the merged units. Final data is output data that is no interim data.
- The new Process has all parameters of the merged units.
- The new Process has all tags of the merged processes, but without duplicates.
- (The Multiplier of a process with multiple Tags is determined by the highest tag-group (*prepare, analyze, model*) present in the Tags.)
- The name of the process of the new Process is equal to the most significant merged process.





Group	Tag (original)	Multiplier	Description	Examples (non-comprehensive)
prepare	selection (Selection)	1	Selection operations retrieve data from the database, [...].	Search, Logical/Arithmetic/Algebraic Condition, Suppress, Filter, Data Retrieval
	transformation (Transformation)	1	Transformation operations transform selected data for scale, orientation, projection, and presentation.	Scale Change, Orientation, Projection Change, Georeferencing, Map Registration, Resampling, raster-vector-conversion
	reclassification (Reclassification)	1	Reclassification operations assign new attribute values to a set of objects based on initial attribute values, geometry or topological relations.	Filtering, Merge Attributes, Generalization
	geometric (Geometric Computation)	1	Geometric Computation operations calculate simple geometric values like area, length, perimeter, and distance, but also more complex calculations like buffer zone, and average distance.	Area, Length, Perimeter, Distance, Buffer Zone, Generalization, Centroid, Dissolve Lines, Thiessen, Merge, Surface Calculation, Image Enhancement/Edge Detection
	visualization (Visualization)	1	Visualization covers any data visualization technique. Used in data explorations, process validation, result visualization, etc.	Block Diagram, Scene Generation, Cross Section, Histogram, Choropleth Map, Error Bars

analyze	neighborhood (Neighborhood Analysis)	2	Neighborhood operations use neighborhood information to calculate slope, diversity, profile, and interpolation.	Slope, Diversity, Interpolation, Proximity, Complex Generalization, Connectivity, Adjacency, Nearest Neighbor, Diffusion Model, Line-of-Sight, Pattern Analysis, Viewshed
	overlay (Map Overlay)	2	Map Overlay combines two or more maps; Both spatial and thematic attribute values can be combined.	Line-on-Polygon, Polygon Overlay, Corridors, Line Intersection, Route Allocation, Map Algebra, Optimal Location
	network (Transitive Closure Analysis)	2	Transitive Closure Analysis [...] calculate shortest path, travel time, and maximum flow. This can be done on a line graph or a polygon network.	Shortest Path, Travel Time, Connectivity, Maximum Flow/Flow Analysis, Drainage Network, Watershed Boundaries, Diffusion Model, Catchment
	statistics (Statistics)	2	Statistics operations calculate statistical characteristics like percentage, distribution, frequency, and correlation of objects of different classes over a specified area.	Percentage, Distribution, Frequency, Correlation, Spatial Statistics, Pattern Analysis, Complexity and Variation Measures
model	model	3	Model describes any black-box module. It is also applicable for scripted processing steps that can not be abstracted by the categories above.	Neural Network,