

2025.09.03

Parsimonious knowledge-based summary of the annotations of a set of elements

- annotations are not independent; they are organized in an ontology
- elements are described using annotations
 - the structure of the ontology is taken into account:
if an element is associated to an annotation,
then it is also associated to all the ancestors
of this annotation
 - each element can be associated (directly or
indirectly) to 0..n annotations
 - likewise, each annotation can describe
(directly or indirectly) 0..n elements
- elements of interest: subset of the set of elements
 - the set of the elements of interest can itself
be also described using annotations.
 - this description can be inferred automatically, e.g. through GSEA
 - **fb-1:** this usually results in numerous annotations,
even after filtering
 - **fb-2:** many of these annotations are interdependent,
so there is some redundancy

Hypothesis: we can exploit the relations between the annotations to identify the most important ones and discard the others.

Greedy algorithm (v2)

Intuition: initially, the summary is empty and the candidates are the annotations over-represented in the set of the elements of interest

- ② identify the annotation(s) that have the best compromise between the information they convey (the more informative, the better) and the number of elements of interest it annotates. That are not annotated by the summary (which favors general so less informative annotations)
- ③ add them to the summary, and discard them from the candidates. Also discard all their descendants from the candidates, as they necessarily annotate a subset of the elements of interest annotated by the annotation(s) we have chosen
- ④ discard all the remaining candidates that do not annotate any element of interest that is not annotated by the summary (X)

- repeat the 3 previous steps until candidates is empty or the remaining candidates do not annotate any element of interest that is not annotated by the summary

\diamond check that guarantees that the 2nd condition would result in candidates = {} at the next step.

(8) prune redundant annotations from the summary (*) i.e.

- annotations from the summary that have an ancestor in the summary and, the ancestor annotates > 1 element of interest that is not annotated by any other annotation from the summary \Rightarrow discard the descendant
- annotations from the summary that have an ancestor in the summary and, all the elements of interest are annotated by > 1 annotation from the summary that is different from the ancestor \Rightarrow discard the ancestor

Δ if > 1 couple (ancestor; descendant) in the summary, the order of pruning probably matters
 \Rightarrow figure out the best strategy
 \diamond try all the combinations, and choose the one that minimizes the size of the summary and maximizes the sum of IC ?

Greedy algorithm (v2)

// initialization

foreach annotation, compute a score that reflects (1) its general relevance and (2) the number of elements of interest it annotates

Todo: try

$$\text{or cumulative IC}(a) = IC(a) \times |\text{annotatedBy}(a) \cap EOI|$$

$$\text{or cumulative Entropy} = H(a) \times |\dots|$$

$$\text{with } p(a) = \frac{|\text{annotatedBy}(a)|}{|\text{Elements}|}$$

$$IC(a) = -\log_2(p(a))$$

$$H(a) = -p(a) \log_2(p(a))$$

candidates = {annotations over represented in the elements of interest}

summary = {}

ItsAnnotatedByCandidates = {el^{ts} of interest annotated by ≥ 1
over-represented annot}

NoAnnotatedBySummary = {}

// incremental aggregation of the summary by pruning
// the candidates

while ($\text{candidates} \neq \{\}$ and $\text{el}^{\text{isAnnotatedByCandidates}} \setminus \text{el}^{\text{isAnnotBySummary}} \neq \{\}$)

compute $\text{scoreMax} = \max(\{\text{score}(c) \mid c \in \text{candidates}\})$
 $cWNS = \{c \in \text{candidates} \mid \text{score}(c) = \text{scoreMax}\}$

remove the annotations from $cWNS$ that have a descendant
in $cWNS$ with the same coverage (x)

remove the descendants of $cWNS$ from candidates

$\text{summary} = \text{summary} \cup cWNS$

$\text{candidates} = \text{candidates} \setminus cWNS$

update $\text{el}^{\text{isAnnotatedByCandidates}}$
 $\text{el}^{\text{isAnnotatedBySummary}}$

$\text{candidates} = \text{candidates} \cap \text{annot}^{-1}(\text{el}^{\text{isAnnotatedByCandidates}} \setminus \text{el}^{\text{isAnnotatedBySummary}})$

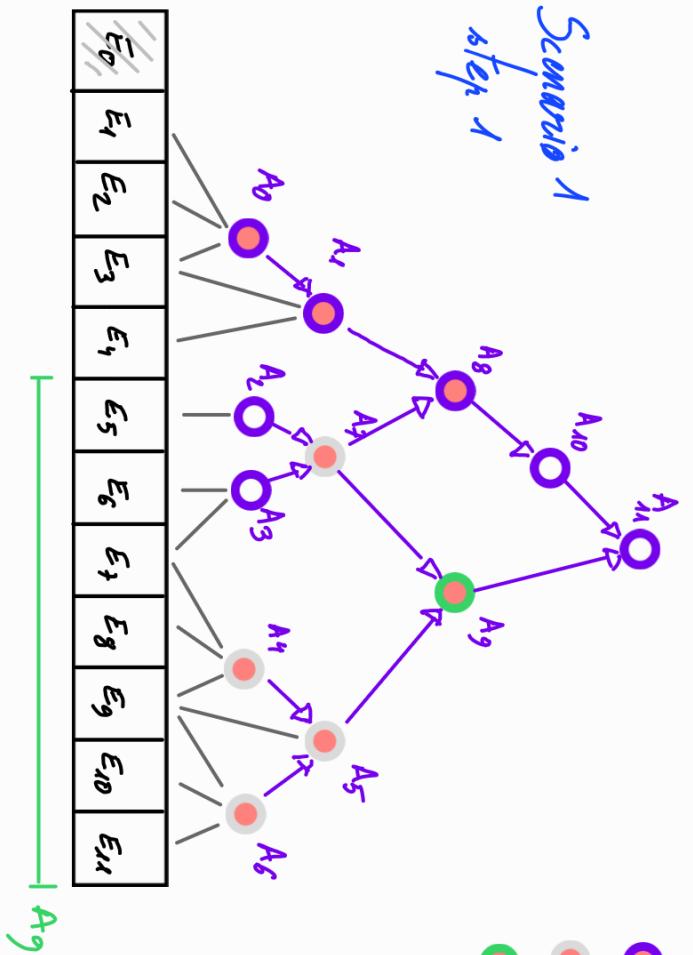
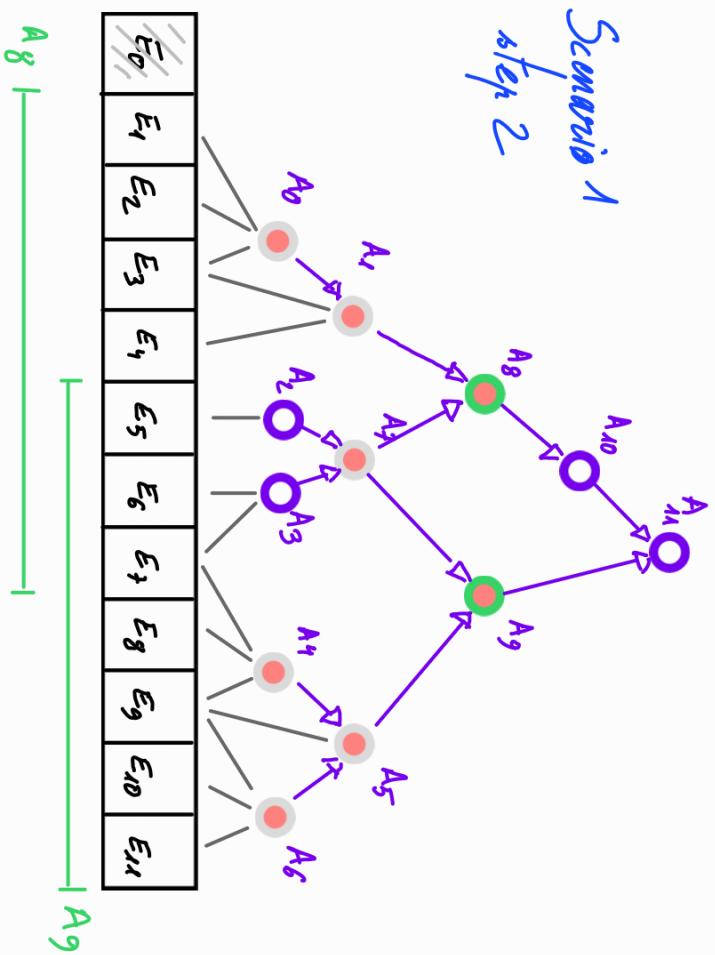
(remove from candidates the annots that are not
associated with any el^{e} not yet annotated by summary)

summarize the summary (remove the redundant annotations, if any) (x)

Greedy algorithm

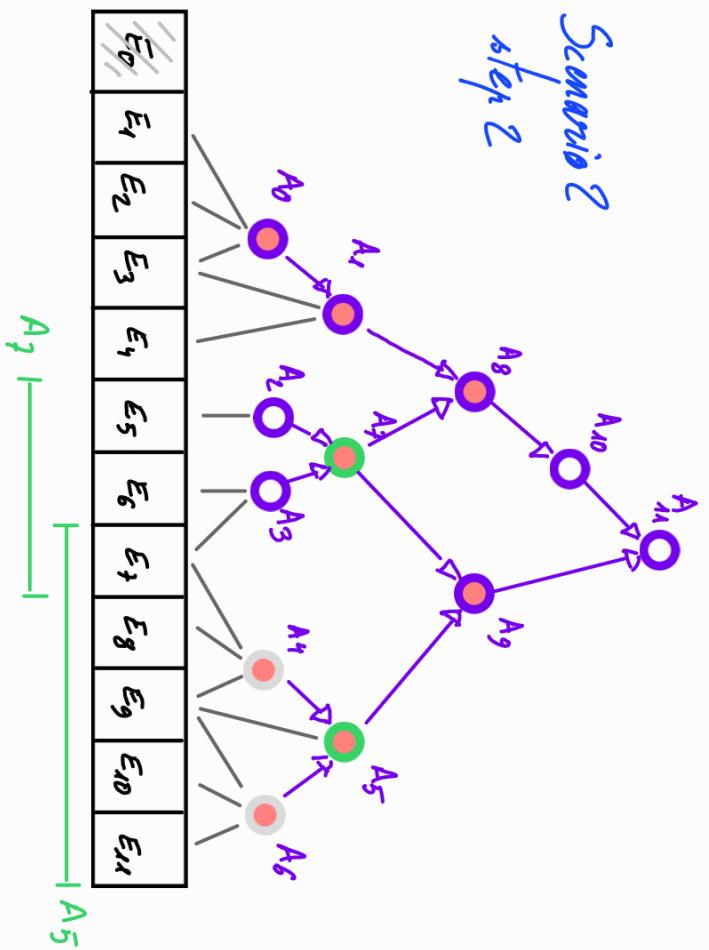
2. ontology-based tricks (x)

- at each step, if multiple candidates with same coverage and one is a superset of the other, keep the most precise (subtree) \Rightarrow informative
- at the end, prune the annotations that are covered by 1+ other annotation of the summary (i.e. all its direct supersets are descendants of the other annotations of the summary) \Rightarrow concision



- As selected and added to the summary
 - no new candidate can reach new elements \Rightarrow end of the algorithm

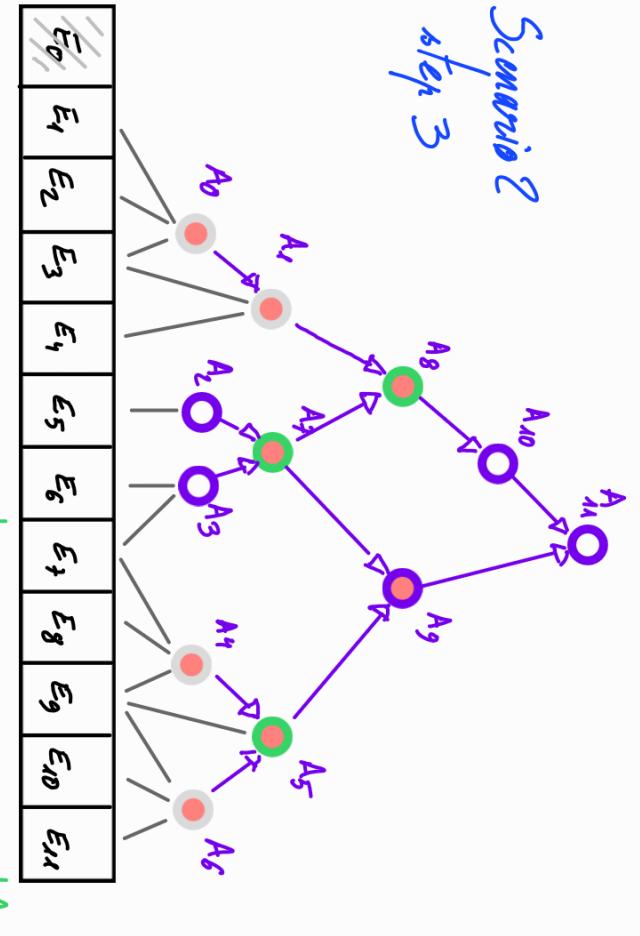
- A_9 selected as first candidate and added to summary
 - its descendants are removed from the pool of candidates



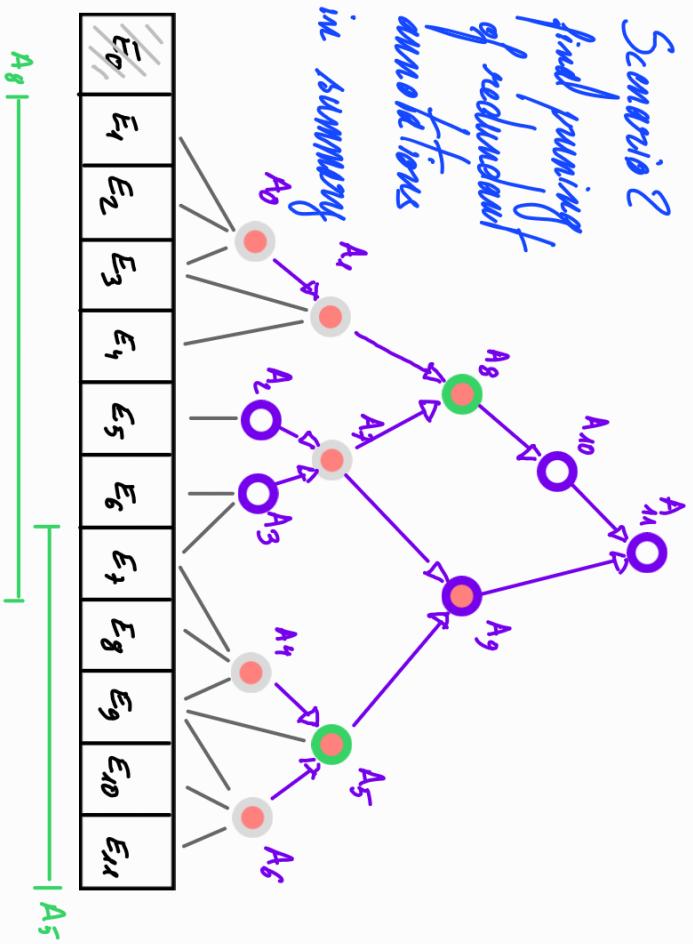
scenario 2: A_7 better candidate than A_9

A_9 still a possible candidate in future steps
but does not annotate new edges so
will never be selected

Scenario 2 step 3

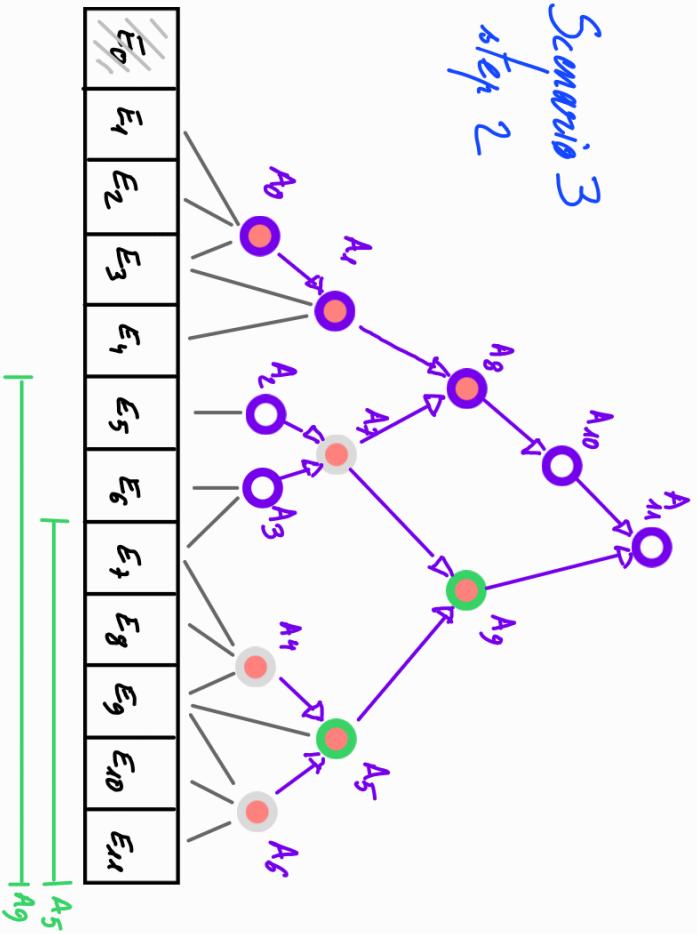


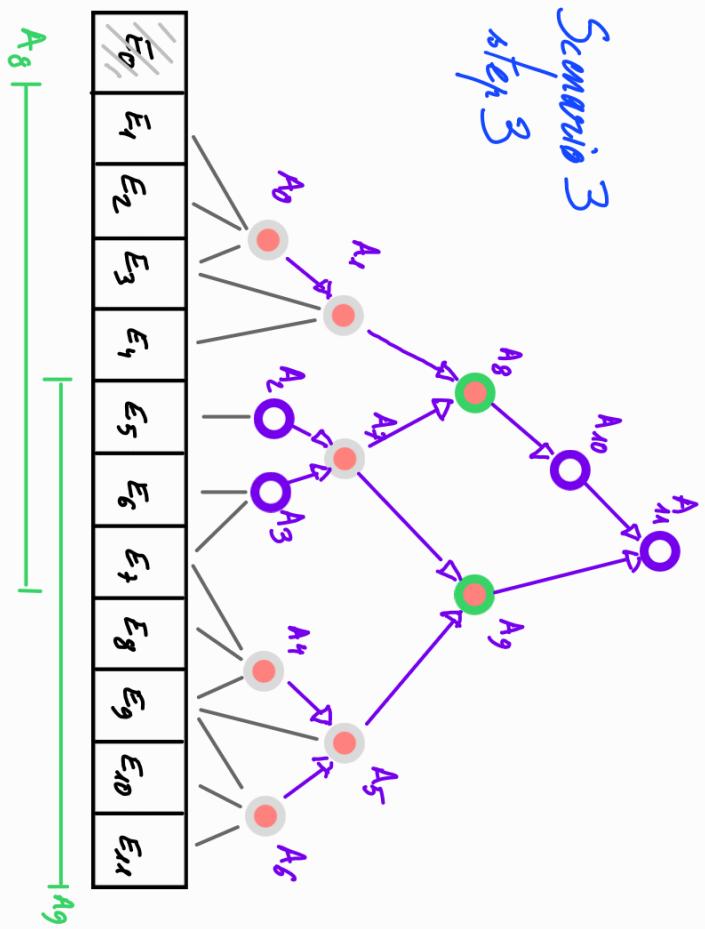
Scenario 2 final pruning of redundant annotations





- contrary to scenario 2 step 2,
 A_9 is a better candidate than A_7





- final warning of redundant annotations in summary:
discard A_5 (or A_9)