

GUARANTEED OPTIMAL COMPOSITIONAL EXPLANATIONS FOR NEURONS

Anonymous authors

Paper under double-blind review

ABSTRACT

While neurons are the basic units of deep neural networks, it is still unclear what they learn and if their knowledge is aligned with that of humans. Compositional explanations aim to answer this question by describing the spatial alignment between neuron activations and concepts through logical rules. These logical descriptions are typically computed via a search over all possible concept combinations. Since computing the **spatial** alignment over the entire state space is computationally infeasible, the literature commonly adopts beam search to restrict the space. However, beam search cannot provide any theoretical guarantees of optimality, and it remains unclear how close current explanations are to the true optimum. In this theoretical paper, we address this gap by introducing the first framework for computing guaranteed optimal compositional explanations. Specifically, we propose: (i) a decomposition that identifies the factors influencing the **spatial** alignment, (ii) a heuristic to estimate the alignment at any stage of the search, and (iii) the first algorithm that can compute optimal compositional explanations within a feasible time. Using this framework, we analyze the differences between optimal and non-optimal explanations **in the most popular settings for compositional explanations, the computer vision domain and Convolutional Neural Networks. In these settings,** we demonstrate that 10–40% of explanations obtained with beam search are suboptimal when overlapping concepts are involved. Finally, we evaluate a beam-search variant guided by our proposed decomposition and heuristic, showing that it matches or improves runtime over prior methods while offering greater flexibility in hyperparameters and computational resources.

1 INTRODUCTION

Compositional explanations (Mu & Andreas, 2020) are a method for **expressing the alignment between the locations of a given neuron activation range and the location of concepts, identified by annotations in datasets.** The key idea is capturing the interaction between low-level (e.g., colors, textures, shapes) and high-level concepts (e.g., objects, entities) via propositional logic formulas able to express the alignment between these complex relationships and neuron activations.

One of the key problems to achieve this goal is that the full search space encompassing all of the possible combinations between concepts cannot typically be exhaustively explored due to its size. As a result, prior work has relied on beam search to identify **the highest spatial alignment within the restricted space**(Mu & Andreas, 2020). Although the resulting explanations are valid, it is unclear whether these explanations are, in fact, optimal (*i.e., whether they identify the specific combination that expressed the highest absolute spatial alignment in the search space, see Figure 1*). While beam search does not guarantee optimality, it might converge to it due to unstudied or unknown properties of the underlying datasets. If not, the explanations produced by beam search may represent only a subset of the alignment structure, offering a partial view of neuron behavior.

The main contribution of this work is to make navigating the state space tractable. To achieve this, we propose a decomposition of the Intersection-over-Union (IoU) metric that reveals a set of fundamental quantities governing alignment quality, and we design a heuristic and a corresponding algorithm that jointly reduce the size of the state space and guide the search process. This approach enables the computation of optimal compositional explanations in feasible time. To the best of our knowledge, this is the first attempt in this research direction. As a first step, we apply our

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107



Figure 1: An example of a case where the optimal algorithm finds a combination of concepts expressing the highest spatial alignment that beam search fails to capture.

method to the computer vision domain and Convolutional Neural Networks, given its prominence in compositional explanation research.

Our contributions are as follows:

- We propose a decomposed Intersection over Union score (dIoU) that identifies fundamental quantities for alignment quality and enables a better characterization of the impact of logical operators on spatial alignment.
- We design a heuristic and a corresponding algorithm that jointly reduce the size of the state space and guide the search process. We show that this algorithm computes **guaranteed optimal** explanations in feasible time and we analyze the differences between optimal and non-optimal explanations.
- We show that part of our proposed heuristic can be used directly to guide beam search with significant gains in flexibility and competitive or better performance than competitors. Specifically, our variant scales more effectively than competitors with respect to explanation length and beam size, and it is less resource-intensive and easier to parallelize.

In Section 2, we give an introduction to the relevant literature in the background and specify our framework for optimal compositional explanations. In Section 3, we analyze our contribution.

2 OPTIMAL COMPOSITIONAL EXPLANATIONS

2.1 BACKGROUND AND RELATED WORK

Neuron explanations aim to understand what individual neurons learn during the training process. Different categories of methods have been proposed to decode different behaviors. Among the most popular in computer vision, we can cite the one that generates samples that capture features recognized by a neuron (Erhan et al., 2009; Olah et al., 2017; Nguyen et al., 2016) and the ones that generate textual descriptions that correlate neuron activations and samples associated with a given concept (Hernandez et al., 2022; Oikarinen & Weng, 2023; 2024; Kalibhat et al., 2023; Ahn et al., 2024; Wang et al., 2022) through foundational models.

Differently from them, compositional explanations are a family of neuron explanations that specifically focus on the **spatial alignment** between a neuron activation range and **and the location of** concepts and express them through propositional logic formulas. The seminal work in this area is Network Dissection (Bau et al., 2017; 2020), which associates each neuron with the single concept that maximizes this alignment. This approach was extended by Mu & Andreas (2020) to associate relationships between multiple concepts, in an attempt to capture a higher degree of polysemantic

behavior (Elhage et al., 2022). Relationships explored in the literature include co-occurrence (Bau et al., 2017), exclusion (Mu & Andreas, 2020), relative position (Harth, 2022), and hierarchy (Massidda & Bacciu, 2023).

Preliminaries and Terminology Let \mathcal{L}^1 be a concept set including properties of interest for a given model and task (e.g., names of colors, objects, categories, shapes). Let $\mathcal{D} = \{x_1, x_2, \dots, x_n\}$ be a dataset, where each input x_i has dimension d and is associated with a set of annotations (possibly empty) identifying the specific locations of concepts in that sample. In this paper, we use the term “location” to refer to the pair (x, j) where x identifies the position of a sample in the probing dataset and j identifies the position of a feature. In vision tasks, these annotations correspond to segmentation masks that group features (i.e., pixels) semantically related, and a “concept” is the semantic label assigned by a human to that annotation. Note that the dataset can differ from the one used to train the model and is assumed to be used solely for neuron analysis. For this reason, \mathcal{D} is commonly referred to as a “probing dataset”. The literature define the *Concept Tensor* $\mathbf{M}^{\mathcal{L}^1} \in \{0, 1\}^{|\mathcal{L}^1| \times |\mathcal{D}| \times d}$ as the binary tensor corresponding to the localization of each concept within the dataset samples, where an element in position $[k, x, j]$ is 1 if the annotations associated with the feature (i.e., pixel) in position j of the sample x include the concept k , and 0 otherwise. From the Concept Tensor, we can extract, for each concept k , the *Concept Matrix* $\mathbf{M}_k \in \{0, 1\}^{|\mathcal{D}| \times d}$. We use the notation $\mathbf{1} \cdot$ to indicate the set of locations equal to 1 in a binary matrix.

Let z be a neuron to be explained in a probed model, and let d be the dimension of its activations. In this paper, following the standard setting in the compositional explanations literature (Mu & Andreas, 2020), z denotes a neuron in a convolutional layer and d is obtained by upscaling its feature map through interpolation. From its activation, we can identify the *Neuron Activation Matrix* $\mathbf{N} \in \{0, 1\}^{|\mathcal{D}| \times d}$ as the binary matrix indicating the locations where the neuron fires within the dataset samples, where an element in position $[x, j]$ is 1 if the activation corresponding to position j in sample x lies within the considered activation range, and 0 otherwise.

Let \mathcal{L}^n be the set of all possible logical formulas of arity at maximum n between concepts in \mathcal{L}^1 . Compositional explanations aim to assign to z the logical combination $L \in \mathcal{L}^n$ of concepts in \mathcal{L}^1 (e.g., ((Cat OR Car) AND White)) that maximizes the alignment between the localization of a given neuron’s activation range $[\tau_1, \tau_2]$ and the localization of the concepts within the probing dataset.

Formally, the algorithm identifies the label $L \in \mathcal{L}^n$ that maximizes the following objective:

$$\arg \max_{L \in \mathcal{L}^n} IoU(L, \mathbf{N}, \mathbf{M}) \quad (1)$$

where the Intersection Over Union (*IoU*) measures the overlap between label annotations and neuron activations, and it is defined as:

$$IoU(L, \mathbf{N}, \mathbf{M}) = \frac{|\mathbf{1}\mathbf{N} \cap \mathbf{1}\theta(\mathbf{M}, L)|}{|\mathbf{1}\mathbf{N} \cup \mathbf{1}\theta(\mathbf{M}, L)|} \quad (2)$$

$|\cdot|$ indicates the cardinality of a set, $\mathbf{1} \cdot$ indicates the set of locations equal to 1 in a binary matrix, and $\theta(\mathbf{M}, L)$ is a function that returns the logical combination of the matrices in \mathbf{M} of the concepts involved in the label L . The label L is typically identified through a search algorithm. To reduce the search space, compositional explanations typically make two assumptions: concepts in the explanation are distinct (**Assumption 1**), and concepts are combined incrementally to form labels (e.g., (((((A \oplus B) \oplus C) . . .) \oplus Z)) (**Assumption 2**), where \oplus indicates a logical connective. Even with these assumptions, \mathcal{L}^n is too large to be explored exhaustively. The total number of combinations is $\sum_{k=1}^n n_o^{k-1} \prod_k (|\mathcal{L}^1| - k)$, where n_o is the number of logic connectives, and each combination requires the comparison of $2d$ values to compute the alignment. In the settings considered by Mu & Andreas (2020), this leads to 2.8×10^{14} operations, rendering both storage and runtime infeasible. To cope with this, prior work adopts vanilla (Mu & Andreas, 2020; Harth, 2022; Massidda & Bacciu, 2023; Makinwa et al., 2022) or informed (La Rosa et al., 2023) beam search with a small beam size. However, beam search does not guarantee optimality, leaving open the question of whether the resulting compositional explanations are the best possible or if better-aligned explanations exist. In the following, we characterize the fundamental quantities that influence alignment and propose both a heuristic and an algorithm that guarantees the optimality of explanations.

2.2 FUNDAMENTAL QUANTITIES

This section introduces our **proposed** decomposed Intersection over Union score (dIoU) and the corresponding fundamental quantities. Alongside the assumptions introduced earlier, we adopt the following:

Assumption 3. *The logic operators connecting concepts are 00-preserving (i.e., they cannot produce a 1 from two zeros).*

This assumption has been implicitly adopted in prior literature, where the commonly used bitwise OR, AND, and AND NOT operators all satisfy the 00-preserving property. **Given the setup described in the previous section**, we propose the definitions of the following quantities.

Definition 2.1. We define the set U of *unique elements* of \mathfrak{D} as the set of locations associated with exactly one annotation, and the set C of *common elements* of \mathfrak{D} as the set of locations associated with multiple annotations.

$$U = \{(x, j) \mid \exists! k \in \mathfrak{L}^1 \text{ s.t. } x \in \mathfrak{D}, j \in [d], \mathbf{M}^{\mathfrak{L}^1}[k, x, j] = 1\}$$

$$C = \{(x, j) \mid \exists k_1, k_2 \in \mathfrak{L}^1 \text{ s.t. } x \in \mathfrak{D}, j \in [d], k_1 \neq k_2, \mathbf{M}^{\mathfrak{L}^1}[k_1, x, j] = 1, \mathbf{M}^{\mathfrak{L}^1}[k_2, x, j] = 1\}$$

In other words, unique elements are features associated with one and only one concept, while common elements are features associated with multiple concepts.

Definition 2.2. Given a neuron activation matrix N , we define the *unique activation* set N^U as the set of locations where the neuron fires on unique elements, and the *common activations* N^C as the set of locations where the neuron fires on common elements.

$$N^U = \{(x, j) \mid (x, j) \in U, N[x, j] = 1\}$$

$$N^C = \{(x, j) \mid (x, j) \in C, N[x, j] = 1\}$$

Definition 2.3. Given a neuron activation matrix N , a concept $k \in \mathfrak{L}^1$, and its **Concept Matrix** M_k , we define the *unique intersection* set I^U as the set of unique elements that are annotated with k and where the neuron fires, and the *common intersection* set I^C as the set of common elements in annotated with k where the neuron fires.

$$I^U(k) = \{(x, j) \mid M_k[x, j] = 1, (x, j) \in N^U\}$$

$$I^C(k) = \{(x, j) \mid M_k[x, j] = 1, (x, j) \in N^C\}$$

Definition 2.4. Given a neuron activation matrix N , a concept $k \in \mathfrak{L}^1$, and its **Concept Matrix** M_k , we define the *unique extras* set E^U as the set of all unique elements in \mathfrak{D} annotated with the concept k and where the neuron does not fire, and the *common extras* set E^C as the set of all common elements in \mathfrak{D} annotated with the concept k where the neuron does not fire.

$$E^U(k) = \{(x, j) \mid M_k[x, j] = 1, (x, j) \in U, (x, j) \notin N^U\}$$

$$E^C(k) = \{(x, j) \mid M_k[x, j] = 1, (x, j) \in C, (x, j) \notin N^C\}$$

Definition 2.3 and Definition 2.4 can be generalized to the case of label $L \in \mathfrak{L}^n$. In this case, the binary label matrix M_L is obtained as the result of the bitwise logic operations induced by the logic operators connecting single concepts $k \in L$.

Definition 2.5. Given a binary neuron activation matrix N and a label $L \in \mathfrak{L}^n$, the decomposed alignment between the annotations associated with L and the neuron activations is defined as:

$$dIoU(N, L, \mathfrak{D}, I^C, I^U, E^U, E^C) = \frac{\sum_{x \in \mathfrak{D}} |I^U(L)_x| + |I^C(L)_x|}{|{}^1N| + \sum_{x \in \mathfrak{D}} |E^U(L)_x| + |E^C(L)_x|} \quad (3)$$

where the subscript indicates the quantity per sample.

Lemma 1. **Given a binary neuron activation matrix N and a label $L \in \mathfrak{L}^n$, the decomposed alignment between the annotations associated with L and the neuron activations is equivalent to the Intersection Over Union if the logical operators involved in L are 00-preserving:**

$$\frac{\sum_{x \in \mathfrak{D}} |I^U(L)_x| + |I^C(L)_x|}{|{}^1N| + \sum_{x \in \mathfrak{D}} |E^U(L)_x| + |E^C(L)_x|} = \frac{|{}^1N \cap {}^1M_L|}{|{}^1N \cup {}^1M_L|} \quad (4)$$

216 *Proof.* See Section A. □

217
218
219 We visualize all the identified quantities and the $dIoU$ score in Table 3 in the Appendix.

220
221 **Observation 1** (Impact of Operators on Quantities). Given the above definitions, we can quantify
222 the impact of the OR, AND, and AND NOT bitwise logic operators connecting two concepts k_1 and
223 k_2 . The OR operator is 1-preserving (i.e., any 1 in either concept remains 1), and thus it preserves
224 the common elements shared by the two concepts and combines (i.e., sums) the non-shared ones
225 and unique elements. The AND operator is 0-preserving (i.e., any 0 in either concept forces 0)
226 and therefore removes all of the unique elements as well as common elements not shared by both
227 concepts. Finally, the AND NOT operator preserves all of the unique elements but removes the
228 common elements shared by both concepts. These operators will be the ones considered in this
229 paper.

230 2.3 HEURISTIC

231
232 This section introduces our proposed heuristic. Given a label $L \in \mathcal{L}^i$ s.t. $i \leq n$, the generic goal
233 of this heuristic is to estimate the label alignment in a faster and less computationally intensive
234 way than directly computing it. Specifically, the proposed heuristic gives an estimation for the
235 following maximum and minimum quantities: (i) given a logic operator \oplus , a label \mathcal{L}^i , and a concept
236 k , the estimate alignment of $L \oplus k$, (ii) alignment reachable starting from L and chaining additional
237 concepts not yet included in L , up to a length of n .

238 2.3.1 ESTIMATE LABEL QUANTITIES

239
240 To facilitate the estimation of the alignment of $L \oplus k$, we introduce the binary *Disjoint Matrix*
241 $D \in \{0, 1\}^{|\mathcal{L}^1| \times |\mathcal{L}^1|}$, which encodes whether two concepts share any annotation overlap and can be
242 computed once per dataset as a preprocessing step. Specifically, for every pair (k_1, k_2) of concepts:

$$243 D[k_1, k_2] = \begin{cases} 1, & \text{if } \forall (x, j) \in |\mathcal{D}| \times d \mid M_{k_1}[x, j] \neq M_{k_2}[x, j], \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

244
245 In addition, we introduce two new quantities, *Space for Common Extras* and *Space for Unique*
246 *Extras*, derived from Definitions 2.1, 2.2 and 2.4, to characterize and constrain the available space
247 for the set of extra elements:

$$248 SE^C = \{(x, j) \mid (x, j) \in C, N[x, j] = 0\} \quad SE^U = \{(x, j) \mid (x, j) \in U, N[x, j] = 0\} \quad (6)$$

249
250 By Assumption 2, we can separate the label into its left (L_{\leftarrow}) and right (L_{\rightarrow}) sides. We can observe
251 that unique elements are always disjoint elements and cannot be shared by definition definition 2.1.
252 Therefore, their value can be computed exactly by using Observation 1.

253
254 **Regarding the common elements**, we can use D to check whether the left and right sides of L are
255 disjoint and estimate alignment differently in the two cases.

256
257 **Disjoint:** If the two sides are disjoint, then common elements follow the same principle as the
258 unique elements since the sides do not share elements. Note that for the AND NOT operator, all the
259 quantities converge to the L_{\leftarrow} , representing a mathematical equivalence and a degenerate case.¹ To
260 force the algorithm to ignore these duplicate and degenerate cases, we manually set their $dIoU$ to 0.

261
262 **Overlap:** If the two sides overlap, we can estimate the exact value of the unique quantities for
263 OR and AND using the same formulas as in the disjoint case. For the AND NOT operator, by
264 Observation 1, the value equals the quantity of the left side. For the common quantities, we can
265 derive them by combining the definitions in Section 2.2 and Observation 1, obtaining:

266
267
268
269 ¹For example, “cat AND NOT dog” is always true, and the AND NOT side does not add meaningful information to the explanation.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286

$$|I_{min}^C(L)_x| = \begin{cases} \max(|I_{min}^C(L_{\leftarrow})_x|, |I^C(L_{\rightarrow})_x|) & \text{if } \oplus = OR \\ \max(|I_{min}^C(L_{\leftarrow})_x| + |I^C(L_{\rightarrow})_x| - |N_x^C|, 0) & \text{if } \oplus = AND \\ \max(|I_{min}^C(L_{\leftarrow})_x| - |I^C(L_{\rightarrow})_x|, 0) & \text{if } \oplus = AND NOT \end{cases} \quad (7)$$

$$|I_{max}^C(L)_x| = \begin{cases} \min(|I_{max}^C(L_{\leftarrow})_x| + |I^C(L_{\rightarrow})_x|, |N_x^C|) & \text{if } \oplus = OR \\ \min(|I_{max}^C(L_{\leftarrow})_x|, |I^C(L_{\rightarrow})_x|) & \text{if } \oplus = AND \\ \min(|I_{max}^C(L_{\leftarrow})_x|, |N_x^C| - |I^C(L_{\rightarrow})_x|) & \text{if } \oplus = AND NOT \end{cases} \quad (8)$$

$$|E_{min}^C(L)_x| = \begin{cases} \max(|E_{min}^C(L_{\leftarrow})_x|, |E^C(L_{\rightarrow})_x|) & \text{if } \oplus = OR \\ \max(|E_{min}^C(L_{\leftarrow})_x| + |E^C(L_{\rightarrow})_x| - |SE_x^C|, 0) & \text{if } \oplus = AND \\ \max(|E_{min}^C(L_{\leftarrow})_x| - |E^C(L_{\rightarrow})_x|, 0) & \text{if } \oplus = AND NOT \end{cases} \quad (9)$$

$$|E_{max}^C(L)_x| = \begin{cases} \min(|E_{max}^C(L_{\leftarrow})_x| + |E^C(L_{\rightarrow})_x|, |SE_x^C|) & \text{if } \oplus = OR \\ \min(|E_{max}^C(L_{\leftarrow})_x|, |E^C(L_{\rightarrow})_x|) & \text{if } \oplus = AND \\ \min(|E_{max}^C(L_{\leftarrow})_x|, |SE_x^C| - |E^C(L_{\rightarrow})_x|) & \text{if } \oplus = AND NOT \end{cases} \quad (10)$$

287 Note that, since L_{\rightarrow} is always an atomic concept (by Assumption 2), its quantities are exact. A
288 detailed derivation is provided in Section E.

289
290 **Aggregated Computation** The above quantities are defined per-sample, requiring $|\mathcal{D}|$
291 comparisons for each computation. This can still be costly for large state spaces. To mitigate
292 this, we introduce a lighter aggregated computation, obtained by summing the values per label.
293 For example, $\sum_{x \in \mathcal{D}} \min(|E_{max}^C(L_{\leftarrow})_x|, |SE_x^C| - |E^C(L_{\rightarrow})_x|)$ can be transformed in
294 $\min(\sum_{x \in \mathcal{D}} |E_{max}^C(L_{\leftarrow})_x|, \sum_{x \in \mathcal{D}} |SE_x^C| - \sum_{x \in \mathcal{D}} |E^C(L_{\rightarrow})_x|)$. The trade-off is precision versus
295 efficiency: the sample version is more accurate but computationally intensive, while the aggregated
296 version can be pre-computed once per label, reducing the cost to a single comparison per quantity
297 at the expense of precision (see Section C for details).

2.3.2 ESTIMATE PATHS

298
299
300 In this section, we estimate the maximum score achievable by concatenating **an arbitrary number of**
301 concepts to a label. If the explanation length were unbounded, this maximum would converge to 1,
302 making the heuristic uninformative. However, compositional explanations are inherently bounded,
303 as long explanations would not aid user understanding; in practice, the maximum length is typically
304 small (e.g., 3). This boundedness allows us to produce tighter estimates. Given a label L , our goal
305 is to estimate the numerator and denominator in Definition 2.5 for all possible paths obtained by
306 concatenating additional concepts through logic operators. To this end, we introduce the *maximum*
307 *and minimum improvement*. For each sample, we extract the *top* $-n$ and *bottom* $-n$ values for each
308 quantity and compute cumulative sums into the *Top* and *Bott* vectors, starting from the largest or
309 smallest values, respectively. For example, $Top_k(I^C)_x$ represents the cumulative sum of the com-
310 mon intersection values of the k concepts with the highest scores for that quantity. These quantities,
311 combined with Observation 1 and the equations in Section 2.3.1, allow us to compute the maximum
312 and minimum factors for OR, AND, and AND NOT exclusive paths (i.e., paths where only one
operator is used from that point onward to concatenate additional concepts):

$$|I_{min}(L)_x| = \begin{cases} \max(|I_{min}^C(L)_x| + |I_{min}^U(L)_x|, Bott_1(I^C)_x + Bott_1(I^U)_x) & \text{OR Path} \\ 0 & \text{AND Path} \\ |I_{min}^U(L)_x| & \text{AND NOT Path} \end{cases} \quad (11)$$

$$|I_{max}(L)_x| = \begin{cases} \min(|I_{max}^C(L)_x| + Top_t(I^C)_x, |N_x^C|) & \text{OR Path} \\ + \min(|I_{max}^U(L)_x| + Top_t(I^U)_x, |N_x^U|) & \text{AND Path} \\ \min(|I_{max}^C(L)_x|, Top_1(I^C)_x) & \text{AND NOT Path} \\ |I_{max}^U(L)_x| + \min(|I_{max}^C(L)_x|, |N_x^C| - Bott_1(I^C)_x) & \end{cases} \quad (12)$$

322
323

(13)

$$\begin{aligned}
|Union_{min}(L)_x| &= |N_x| + \begin{cases} \max(|E_{min}^C(L)_x| + |E_{min}^U(L)_x|, \\ Bott_1(E^C)_x + Bott_1(E^U)_x) & \text{OR Path} \\ 0 & \text{AND Path} \\ |E_{min}^U(L)_x| & \text{AND NOT Path} \end{cases} & (14) \\
|Union_{max}(L)_x| &= |N_x| + \begin{cases} \min(|E_{max}^C(L)_x| + Top_t(E^C)_x, |SE_x^C|) \\ + \min(|E_{max}^U(L)_x| + Top_t(E^U)_x, |SE_x^U|) & \text{OR Path} \\ \min(|E_{max}^C(L)_x|, Top_1(E^C)_x) & \text{AND Path} \\ |E_{max}^C(L)_x| + \min(|E_{max}^U(L)_x|, \\ |SE_x^C| - Bott_1(E^C)_x) & \text{AND NOT Path} \end{cases} & (15)
\end{aligned}$$

where t denotes the difference between the maximum length and the length of the label L . To estimate the values for paths involving multiple operators, we take the maximum and minimum of each quantity across the operators considered (see Section E.3 for a discussion about explicitly modeling every possible combination). Among the paths, we also include the *final path*, computed by Definition 2.5, to denote the case where the label is not further expanded.

These estimates are finally used to compute the maximum and the minimum dIoU:

$$dIoU_{max} = \frac{\sum_{x \in \mathcal{D}} |I_{max}(L)_x|}{\sum_{x \in \mathcal{D}} |Union_{min}(L)_x|} \quad dIoU_{min} = \frac{\sum_{x \in \mathcal{D}} |I_{min}(L)_x|}{\sum_{x \in \mathcal{D}} |Union_{max}(L)_x|} \quad (16)$$

Aggregated Computation As in Section 2.3.1, the aggregated computation estimates the quantities more efficiently by operating on aggregate values (i.e., sums) per label instead of computing them on a per-sample basis. For example, rather than using $|Union_{max}^{sample}(L)_x| = |N_x| + \min(|E_{max}^C(L)_x|, Top_1(E^C)_x)$ the aggregated formulation uses $|Union_{max}^{aggr}(L)| = N + \min(\sum_{x \in \mathcal{D}} |E_{max}^C(L)_x|, Top_1^A(E^C))$, where Top^A is computed concept-wise.

2.4 OPTIMAL ALGORITHM

The optimal algorithm is a best-first search guided by our proposed heuristic. We provide a textual overview of the main steps below and a more detailed discussion and pseudocode in Section C.

1. Compute the exact quantities (Section 2.2) for every concept in the dataset.
2. For each concept, compute $dIoU_{max}$ and $dIoU_{min}$ for all possible paths starting from it, using the heuristic aggregated computation.
3. Initialize the frontier with all paths whose estimated $dIoU_{max}$ is greater than the global maximum of the minimum estimates.
4. Iteratively pop nodes from the frontier, starting from those with the highest estimated $dIoU_{max}$.
 - (a) If the estimate is aggregated, refine it by computing the sample-based estimate and reinsert the node. Otherwise, proceed to the next step.
 - (b) If the node path is not a final one (i.e., can still be extended), expand its label by concatenating every possible concept with every allowed connective. For each new label, compute $dIoU_{max}$ and $dIoU_{min}$ via the heuristic aggregated computation and insert them into the frontier.
 - (c) If the node path is a final one, compute its exact IoU. During this step, the algorithm stores the intermediate quantities of the sub-labels composing the label. This information is then backpropagated to the frontier: nodes that share sub-labels with the evaluated node update their estimates using these exact quantities.
 - (d) Continue until the frontier is empty. During the process, whenever a new maximum of the $dIoU_{min}$ estimates is found, prune the frontier by removing all nodes whose $dIoU_{max}$ falls below this threshold.

Additionally, to reduce redundant computation, the algorithm incorporates a limited set of logical equivalence rules, which are applied before and during the expansion phase, and maintains a buffer to cache recently explored nodes associated with the same estimated $dIoU_{max}$. We discuss these

Table 1: Average number of visited, expanded, and estimated nodes, along with runtime per unit (in minutes), by the optimal algorithm, a beam search guided by our heuristic, and two alternative beam search algorithms.

Algorithm	Optimal	Visited	Expanded	Estimated	Time (sec)
Low Complexity					
Optimal (our)	✓	1	100	778	0.08
Beam + Our H.	✗	1	11	405	0.09
MMESH Beam	✗	135	14.94	716	10.01
Vanilla Beam	✗	716	14.94	-	2.77
Intermediate Complexity					
Optimal (our)	✓	1	2885	1.76×10^6	69.27
Beam + Our H.	✗	1.94	11	24730	9.46
MMESH Beam	✗	41.94	15	37978	37.31
Vanilla Beam	✗	37979	15	-	450
High Complexity					
Optimal (our)	✓	47.36	3.54×10^5	1.28×10^8	5768
Beam + Our H.	✗	5.82	11	28947	139
MMESH Beam	✗	44.98	13.47	53774 ± 2	106
Vanilla Beam	✗	53775	13.47	-	5929

details and justify the design choices in Section C. **Because the maximum $dIoU$ of each path is an overestimation, and since the algorithm is designed to visit all nodes whose $dIoU$ exceeds that of the current best explanation, the algorithm is guaranteed to return the most aligned explanation. In other words, the returned explanation is always the optimal one (proof in Section F).**

3 ANALYSIS

3.1 FEASIBILITY OF OPTIMALITY AND HEURISTIC-GUIDED BEAM SEARCH

This section evaluates the feasibility of the proposed optimal algorithm and the effectiveness of the heuristic when used to guide beam search. We consider three scenarios that vary in annotation complexity: **low, moderate, and high**. The low-complexity setting (Cityscapes (Cordts et al., 2016)) includes a small number of concepts (25), all of which are disjoint. The moderate-complexity setting (the extended version of Ade20K (Zhou et al., 2017b)) includes a much larger number of concepts (847) but with no overlapping annotations. Finally, the high complexity setting (Brodén (Bau et al., 2017)) involves frequent overlaps combined with a large number of concepts (1198). As reference baselines, we include the MMESH-guided beam search (La Rosa et al., 2023) and the vanilla beam search (Mu & Andreas, 2020). The beam search variant that uses our heuristic (*Beam + Our H.*), replaces MMESH with label-quantity estimation (see Section D for details). For each setting, we report the average number (over 50 units) of visited nodes (i.e., those for which the exact IoU is computed), expanded nodes, estimated nodes, and the computation time per unit (std. dev. can be found in Section B). Following Mu & Andreas (2020), we extract 50 random units of the last convolutional layer in a ResNet (He et al., 2016) model trained on Places365 (Zhou et al., 2017a) and use the highest activations (top 0.005 percentile) as activation ranges.

Table 1 shows that the optimal algorithm consistently finds the optimal solution within feasible runtimes across all scenarios. As expected, informed beam search methods are faster, since they explore a much smaller portion of the state space (i.e., estimated nodes). However, the performance of the optimal algorithm remains comparable to the vanilla beam search of Mu & Andreas (2020), even in the most complex settings. Importantly, the number of expanded states is a small fraction of the overall state space (less than 0.1%) and is significantly smaller than the number of estimated states. This property is crucial for refining heuristic estimates without compromising runtime efficiency (see Section C).

Table 2: Percentage of changed explanations of explanations falling into categories 1, 2, and 3 for different models.

Dataset	Diff	Beam IoU	Optimal IoU	Cat 1 (%)	Cat 2 (%)	Cat 3 (%)
ResNet	8%	0.077	0.083	85	4	11
AlexNet	22%	0.045	0.047	93	3	4
DenseNet	39%	0.039	0.041	73	0	27

More notably, beam search guided by our heuristic outperforms all baselines across all settings in terms of visited states while achieving comparable or better runtimes. **Note that both MMESH and our beam variant converge to the same solutions as those found by the extensive search (Vanilla Beam), since they use admissible heuristics.** Compared to MMESH, our approach offers several improvements. First, MMESH relies on annotations during beam expansion, which requires annotations to be kept in memory and GPU resources, thus limiting scalability and parallelization. Conversely, our variant uses annotations only when visiting states, allowing them to be loaded from disk on demand. This design, combined with its higher efficiency, removes the need to store annotations in memory and enables easier parallelization in low-resource scenarios. Our beam variant also scales efficiently with changes in hyperparameters. For example, when varying explanation length (3,5,10, and 20) and beam size (5, 10, and 20) in the moderate settings, the runtime of our variant is stable between 0.16 and 0.18 min/unit. By contrast, MMESH slows down progressively with both explanation length [0.62, 1.28, 4.55, > 240] and beam size [0.62, 1.24, > 240], starting from 0.62 for 3-concepts explanations and beam size fixed to 5 to 4.5 min/units (see Table 6 for a table of results). For explanation lengths or beam sizes higher than 10, MMESH becomes infeasible, despite the limited complexity of the considered settings.

3.2 EXPLANATION ANALYSIS

This section addresses the question we originally posed about beam search-based algorithms: are the explanations they compute optimal? In general, the answer is no. Although beam search often identifies valid explanations, our analysis (Table 2) shows that **in the High Complexity settings** between 10% and 40% of them differ from the optimal ones across several models studied in previous literature (Mu & Andreas, 2020). We classify these differences into three categories: (1) explanations differ in both concepts and IoU, (2) explanations involve the same concepts but differ in how they are connected, resulting in different IoU, and (3) explanations share the same IoU but differ in the way the concepts are connected.

The first category is the most prominent and **often** involves explanations connected by AND and AND NOT operators, suggesting that beam search struggles to express explanations that describe units specialized in recognizing complex scenarios. This discrepancy does not imply that the explanations produced by beam search are incorrect; they still capture meaningful alignments between the unit and the concepts expressed in the formulas. However, they often fall short of reflecting the highest degree of alignment that the unit actually exhibits. The second category includes cases where the optimal explanations are more precise and more aligned (e.g., from ((table OR sink) AND white-c) (IoU=0.036) to ((white-c AND table) OR sink) (IoU=0.040)). These differences are small and, in this case, we can consider the beam search explanations quite close and faithful to the optimal explanations.

Finally, the third category includes cases where the semantics of the explanations changes. For example, consider the explanation ((ball_pit-s OR flower) AND NOT dining_room-s). At first glance, one might interpret this unit as being specialized in recognizing *ball_pit* not located in dining rooms. However, inspecting the dataset reveals that these concepts never co-occur (i.e., they are disjoint), whereas there are instances including both flowers and dining rooms. Thus, part of the explanation is effectively “unverified”. In contrast, the optimal algorithm correctly identifies the alignment as ((flower AND NOT dining_room-s) OR ball_pit-s) and exposes a key limitation of beam search: it cannot backtrack on earlier decisions, and the search may compensate for errors by producing explanations that rely on unverified scenarios (further details and examples for 10 units per model are provided in Section G).

3.3 ALGORITHM ANALYSIS: INSIGHTS AND LIMITATIONS

This section briefly discusses the key design choices behind the optimal algorithm and provides insights into its limitations.

Overall, the beam search guided by our heuristic represents a safe choice when the reduced time for computation is a priority. Conversely, the optimal algorithm represents a first promising step towards the guarantee of optimality on neural explanations and a better choice when the optimality is considered more important than the time of execution. Regarding design choices, Section C includes a detailed discussion of them. In summary, all the quantities and steps proposed in this algorithm are necessary for the search to be tractable in high complexity scenarios. For example, we use aggregated computations to decide which nodes enter the frontier and sample-based ones to parse it, since estimated nodes far outnumber those actually expanded. Importantly, nodes added to the frontier may still be pruned if their estimated IoU falls below the best found so far. Back-propagation and minimum estimation are crucial in this setting, since they reduce redundant nodes and prevent exhaustive exploration of overestimated candidates. By empirically analyzing the algorithm’s execution traces, we identified the following insights and areas of improvement for future research:

Convergence Towards Breath-first search We noted that the state space is explored similarly to breadth-first search, since the top vector dominates the search and roughly overestimates the maximum improvement. In theory, this reliance could make it infeasible for long explanations (although shorter explanations are generally preferred). To mitigate this problem, future research could explore vectors that are (1) specific to a label (which is costly) or (2) novel and better representations of the maximum improvement.

Unmeaningful Units Our algorithm could become much slower when either a unit is not interpretable (the IoU < 0.04 (Bau et al., 2020)) or it is unspecialized (using default rules) and the probing dataset is of high complexity. In these extreme cases, the space to be explored is very large since there is no clear alignment and the combinations of concepts are all similar. One possible solution is to run beam search and then refine the units deemed interpretable for optimality. Alternatively, one could automatically switch from the optimal algorithm to beam search once the frontier grows too large, and use the beam search output to initialize and guide the optimal search.

4 CONCLUSION

This paper presents the first attempt to guarantee optimality in compositional explanations. Specifically, we identified and formalized the fundamental quantities governing spatial alignment, proposed a heuristic to estimate the potential alignment from any label in the search space, and developed an algorithm capable of computing optimal compositional explanations within feasible runtimes. We further demonstrated that our heuristic can also improve existing beam search-based approaches for non-optimal explanations. Since our method does not rely on spatial information, the proposed heuristic is broadly applicable across domains. Moreover, our theoretical contribution may extend beyond neural explanations, as bitwise computations are of interest in areas such as semantic segmentation and communication. Finally, we call for further research on refining this heuristic and on designing new algorithms for computing alternative forms of compositional explanations.

5 REPRODUCIBILITY STATEMENT

To ensure reproducibility, we will release the code upon acceptance. Additionally, we describe all the assumptions of this work in Section 2.1 and Section 2.2. The full pseudo-code for both the optimal algorithm and the beam search guided by our heuristic is provided in Section C and Section D. Proofs and derivations of the estimations are given in Section 2.2, Section 2.3.1, and Section 2.3.2 and Sections A and E. We detail the rationale behind design choices in Section C and provide example outputs of our algorithm in Section G. These examples can serve as gold references when re-implementing the algorithm. Finally, we describe the datasets and additional setup information in Section B.

REFERENCES

- 540
541
542 Yong Hyun Ahn, Hyeon Bae Kim, and Seong Tae Kim. Www: a unified framework for explaining
543 what where and why of neural networks by interpretation of neuron concepts. In *Proceedings of*
544 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10968–10977, 2024.
- 545 David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection:
546 Quantifying interpretability of deep visual representations. In *Computer Vision and Pattern*
547 *Recognition*, 2017.
- 548
549 David Bau, Jun-Yan Zhu, Hendrik Strobelt, Agata Lapedriza, Bolei Zhou, and Antonio Torralba.
550 Understanding the role of individual units in a deep neural network. *Proceedings of the National*
551 *Academy of Sciences*, 2020. ISSN 0027-8424. doi: 10.1073/pnas.1907375117. URL [https://](https://www.pnas.org/content/early/2020/08/31/1907375117)
552 www.pnas.org/content/early/2020/08/31/1907375117.
- 553 Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Con-
554 nerly, Nick Turner, Cem Anil, Carson Denison, Amanda Aspell, Robert Lasenby, Yifan Wu,
555 Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex
556 Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter,
557 Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language
558 models with dictionary learning. *Transformer Circuits Thread*, 2023. [https://transformer-](https://transformer-circuits.pub/2023/monosemantic-features/index.html)
559 [circuits.pub/2023/monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 560 Kirill Bykov, Laura Kopf, Shinichi Nakajima, Marius Kloft, and Marina Höhne. La-
561 beling neural representations with inverse recognition. In A. Oh, T. Naumann,
562 A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural In-*
563 *formation Processing Systems*, volume 36, pp. 24804–24828. Curran Associates, Inc.,
564 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/](https://proceedings.neurips.cc/paper_files/paper/2023/file/4e52bbb99690d1e05c7ef7b4c8b3569a-Paper-Conference.pdf)
565 [file/4e52bbb99690d1e05c7ef7b4c8b3569a-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/4e52bbb99690d1e05c7ef7b4c8b3569a-Paper-Conference.pdf).
- 566
567 Stephen Casper, Tilman Rauker, Anson Ho, and Dylan Hadfield-Menell. Sok: Toward transparent
568 AI: A survey on interpreting the inner structures of deep neural networks. In *First IEEE Con-*
569 *ference on Secure and Trustworthy Machine Learning*, 2023. URL [https://openreview.](https://openreview.net/forum?id=8C5zt-0Utdn)
570 [net/forum?id=8C5zt-0Utdn](https://openreview.net/forum?id=8C5zt-0Utdn).
- 571 Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo
572 Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic
573 urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern*
574 *Recognition (CVPR)*, 2016.
- 575
576 Maximilian Dreyer, Erblina Purlaku, Johanna Vielhaben, Wojciech Samek, and Sebastian La-
577 puschkin. Pure: Turning polysemantic neurons into pure features by identifying relevant circuits.
578 *arXiv preprint arXiv:2404.06453*, 2024.
- 579 Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec,
580 Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish,
581 Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superpo-
582 sition. *Transformer Circuits Thread*, 2022.
- 583
584 Dumitru Erhan, Yoshua Bengio, Aaron Courville, and Pascal Vincent. Visualizing higher-layer
585 features of a deep network. 2009.
- 586
587 Leo Gao, Tom Dupre la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever,
588 Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. In *The Thirteenth Inter-*
589 *national Conference on Learning Representations*, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=tcsZt9ZNKD)
590 [forum?id=tcsZt9ZNKD](https://openreview.net/forum?id=tcsZt9ZNKD).
- 591 Leilani H. Gilpin, David Bau, Ben Z. Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal.
592 Explaining explanations: An overview of interpretability of machine learning. In *2018 IEEE 5th*
593 *International Conference on Data Science and Advanced Analytics (DSAA)*. IEEE, oct 2018. doi:
10.1109/dsaa.2018.00018.

- 594 Rafael Harth. *Understanding Individual Neurons of ResNet Through Improved Compositional For-*
595 *mulas*, pp. 283–294. Springer International Publishing, 2022. ISBN 9783031092824. doi:
596 10.1007/978-3-031-09282-4_24.
597
- 598 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
599 nition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.
600 770–778. IEEE, 2016.
601
- 602 Evan Hernandez, Sarah Schwettmann, David Bau, Teona Bagashvili, Antonio Torralba, and Jacob
603 Andreas. Natural language descriptions of deep features. In *International Conference on Learning*
604 *Representations*, 2022. URL <https://openreview.net/forum?id=NudBMY-tzDr>.
- 605 Robin Hesse, Jonas Fischer, Simone Schaub-Meyer, and Stefan Roth. Disentangling polysemantic
606 channels in convolutional neural networks. In *The First Workshop on Mechanistic Interpretability*
607 *for Vision*, 2025.
608
- 609 Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected
610 convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
611 *Recognition (CVPR)*, July 2017.
- 612 Neha Kalibhat, Shweta Bhardwaj, C Bayan Bruss, Hamed Firooz, Maziar Sanjabi, and Soheil Feizi.
613 Identifying interpretable subspaces in image representations. In *International Conference on*
614 *Machine Learning*, pp. 15623–15638. PMLR, 2023.
615
- 616 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
617 convolutional neural networks. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger
618 (eds.), *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc.,
619 2012. URL [https://proceedings.neurips.cc/paper_files/paper/2012/](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf)
620 [file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf).
- 621 Biagio La Rosa, Leilani H. Gilpin, and Roberto Capobianco. Towards a fuller understanding
622 of neurons with clustered compositional explanations. In *Thirty-seventh Conference on Neu-*
623 *ral Information Processing Systems*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=51PLYhMFwz)
624 [51PLYhMFwz](https://openreview.net/forum?id=51PLYhMFwz).
- 625 Sayo M. Makinwa, Biagio La Rosa, and Roberto Capobianco. Detection accuracy for evaluating
626 compositional explanations of units. In *AIxIA 2021 - Advances in Artificial Intelligence*, pp. 550–
627 563. Springer International Publishing, 2022. doi: 10.1007/978-3-031-08421-8_38.
628
- 629 Riccardo Massidda and Davide Bacciu. Knowledge-driven interpretation of convolutional neural
630 networks. In *Machine Learning and Knowledge Discovery in Databases*, pp. 356–371. Springer
631 International Publishing, 2023. doi: 10.1007/978-3-031-26387-3_22.
632
- 633 Jesse Mu and Jacob Andreas. Compositional explanations of neurons. In *Proceedings of the 34th*
634 *International Conference on Neural Information Processing Systems*, NIPS ’20, Red Hook, NY,
635 USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- 636 Anh Nguyen, Jason Yosinski, and Jeff Clune. Multifaceted feature visualization: Uncovering the
637 different types of features learned by each neuron in deep neural networks. *Visualization for Deep*
638 *Learning workshop, ICML 2016*, February 2016.
639
- 640 Tuomas Oikarinen and Tsui-Wei Weng. CLIP-dissect: Automatic description of neuron representa-
641 tions in deep vision networks. In *International Conference on Learning Representations*, 2023.
642 URL <https://openreview.net/forum?id=iPWiwWHc1V>.
- 643 Tuomas Oikarinen and Tsui-Wei Weng. Linear explanations for individual neurons. In *Forty-first*
644 *International Conference on Machine Learning*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=WIbntm28cM)
645 [forum?id=WIbntm28cM](https://openreview.net/forum?id=WIbntm28cM).
646
- 647 Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2(11), nov
2017. doi: 10.23915/distill.00007.

648 Laura O’Mahony, Vincent Andrearczyk, Henning Müller, and Mara Graziani. Disentangling neuron
649 representations with concept vectors. In *Proceedings of the IEEE/CVF Conference on Computer
650 Vision and Pattern Recognition*, pp. 3769–3774, 2023.

651
652 Laura O’Mahony, Nikola S. Nikolov, and David J.P. O’Sullivan. Towards utilising a range of neural
653 activations for comprehending representational associations. In *2025 IEEE/CVF Winter Confer-
654 ence on Applications of Computer Vision (WACV)*, pp. 2495–2506. IEEE, February 2025. doi:
655 10.1109/wacv61041.2025.00248.

656 Anvita A. Srinivas, Tuomas Oikarinen, Divyansh Srivastava, Wei-Hung Weng, and Tsui-Wei Weng.
657 Sand: Enhancing open-set neuron descriptions through spatial awareness. In *2025 IEEE/CVF
658 Winter Conference on Applications of Computer Vision (WACV)*, pp. 2993–3002. IEEE, February
659 2025. doi: 10.1109/wacv61041.2025.00296.

660 Andong Wang, Wei-Ning Lee, and Xiaojuan Qi. Hint: Hierarchical neuron concept explainer. In
661 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
662 pp. 10254–10264, June 2022.

663
664 Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2.
665 <https://github.com/facebookresearch/detectron2>, 2019.

666
667 Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 mil-
668 lion image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine
669 Intelligence*, 2017a.

670 Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene
671 parsing through ade20k dataset. In *2017 IEEE Conference on Computer Vision and Pattern
672 Recognition (CVPR)*. IEEE, July 2017b. doi: 10.1109/cvpr.2017.544.

673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

Table 3: Visualization of identified quantities for a sample x . In pink the unique extras. In yellow the common extras. In green the unique intersection. In cyan the common intersection.

	Vector						$dIoU$
$N(x)$	1	1	1	0	0	0	
$M_{c_1}[x]$	1	1	0	0	1	1	2/5
$M_{c_2}[x]$	1	1	0	1	0	0	2/4
$M_{c_3}[x]$	1	0	1	0	1	1	2/5

A PROOF EQUIVALENCE DECOMPOSED ALIGNMENT-IOU SCORE

Proof. Let M_L be the binary label tensor representing the result of the bitwise logic operations induced by the logic operators connecting single concepts $k \in L$. The Intersection over Union (IoU) score is then defined as:

$$IoU(N, L, \mathfrak{D}, M_L) = \frac{|{}^1N \cap {}^1\theta(\mathbf{M}, L)|}{|{}^1N \cup {}^1\theta(\mathbf{M}, L)|} \quad (17)$$

We first observe that $|{}^1N \cap {}^1M_L|$ is necessarily equal to $\sum_{x \in \mathfrak{D}} |I^U(L)_x| + |I^C(L)_x|$. Indeed, all elements in $|{}^1N \cap {}^1M_L|$ are associated with both the neuron and the label L . Consider a generic element $j \in {}^1N \cap {}^1M_L$. Since the logic operators are 00-preserving (Assumption 2), this element must be associated with at least one annotation of one of the concepts in L . But if the element is associated with at least a concept, then it is either associated with exactly one, and thus $j \in I^U$, or with multiple concepts, and thus $j \in I^C$. This also holds in the case of concepts connected by the AND NOT operator, because by design this operator must be chained to a positive concept, and AND NOT is 00-preserving.

Regarding the denominator, we can note that

$$|{}^1N \cup {}^1M_L| = |{}^1N| + |{}^1M_L| - |{}^1N \cap {}^1M_L| \quad (18)$$

Observe that for all $j \in {}^1N \cap {}^1M_L$, the contribution $|{}^1M_L| - |{}^1N \cap {}^1M_L| = 0$, since these elements are counted in both sets. Hence, $|{}^1M_L| - |{}^1N \cap {}^1M_L|$ represents the number of elements that are labeled as 1 in M_L but for which the neuron does not fire. This definition coincides with Definition 2.4. Therefore, we can rewrite the denominator as

$$|{}^1N \cup {}^1M_L| = |{}^1N| + \sum_{x \in \mathfrak{D}} |E(L)_x|. \quad (19)$$

Similarly to the numerator case, because the logic operators are 00-preserving, every element in $E(L)_x$ is either included in $E^U(L)_x$ or in $E^C(L)_x$. Thus,

$$|{}^1N \cup {}^1M_L| = |{}^1N| + \sum_{x \in \mathfrak{D}} (|E^U(L)_x| + |E^C(L)_x|). \quad (20)$$

and thus

$$dIoU(N, L, \mathfrak{D}, I^C, I^U, E^U, E^C) = IoU(N, L, \mathfrak{D}, \mathbf{M}). \quad (21)$$

□

B COMPLETE RESULTS

Table 4 reports the averages and standard deviations of the results presented in Table 1. Table 6 shows differences in scalability between beam search guided by our proposed heuristic and beam search guided by MMESH. Both tables are computed over 50 units randomly extracted from the penultimate layer of a ResNet18 model trained on Places365, consistent with prior work on compositional explanations (Mu & Andreas, 2020; Harth, 2022; Makinwa et al., 2022; La Rosa et al.,

Table 4: Average number of visited, expanded, and estimated nodes, along with runtime per unit (in minutes) and standard deviation. Statistics are computed across three settings and datasets: low complexity (Cityscapes), intermediate complexity (Ade20K Full), and high complexity (Brodén). The table reports results for our two proposed algorithms, the beam search powered by MMESH, and the vanilla beam search.

Algorithm	Optimal	Visited	Expanded	Estimated	Time (sec)
Low Complexity					
Optimal (our)	✓	1	100 ± 34	778 ± 221	0.08 ± 0.01
Beam + Our H.	✗	1 ± 0.30	11 ± 0.42	405 ± 25	0.09 ± 0.01
MMESH Beam	✗	135 ± 3	14.94 ± 0.42	716 ± 17	10.01 ± 0.30
Vanilla Beam	✗	716 ± 17	14.94 ± 0.42	-	2.77 ± 0.16
Intermediate Complexity					
Optimal (our)	✓	1	2885 ± 2583	1.76 × 10 ⁶	69.27 ± 40.60
Beam + Our H.	✗	1.94 ± 0.24	11	24730 ± 1090	9.46 ± 0.38
MMESH Beam	✗	41.94 ± 156	15	37978 ± 1.79	37.31 ± 2.87
Vanilla Beam	✗	37979 ± 1.74	15	-	450 ± 10
High Complexity					
Optimal (our)	✓	47.36 ± 99.27	3.54 × 10 ⁵	1.28 × 10 ⁸	5768 ± 1297
Beam + Our H.	✗	5.82 ± 6.47	11	28947 ± 2261	139 ± 23
MMESH Beam	✗	44.98 ± 69.11	13.47 ± 4.54	53774 ± 2	106 ± 18
Vanilla Beam	✗	53775 ± 1	13.47 ± 4.54	-	5929 ± 548

Table 5: Percentage of changed explanations of explanations falling into categories 1, 2, and 3 for different models along with average IoU and standard deviation of these explanations.

Dataset	Diff	Beam IoU	Optimal IoU	Cat 1 (%)	Cat 2 (%)	Cat 3 (%)
ResNet	8%	0.077 ± 0.044	0.083 ± 0.050	85	4	11
AlexNet	22%	0.045 ± 0.023	0.047 ± 0.024	93	3	4
DenseNet	39%	0.039 ± 0.014	0.041 ± 0.015	73	0	27

2023). We follow Bau et al. (2020) and Mu & Andreas (2020) and use the highest activations corresponding to the top 0.005 percentile across the probing dataset as the activation range and fix the maximum explanation length to 3. As probing datasets, we used Cityscapes (Cordts et al., 2016) (accessible at: <https://www.cityscapes-dataset.com/> under MIT License) for the low complexity settings, Ade20kFull (Zhou et al., 2017b) (accessible via the Detectron2 (Wu et al., 2019) framework) for the moderate settings, and Broden (Bau et al., 2017) (accessible at <https://github.com/CSAILVision/NetDissect> under MIT license) for the highest complexity settings. All results were computed on a workstation equipped with an NVIDIA GTX 3090 GPU, without parallelization, in order to avoid timing overhead.

C OPTIMAL ALGORITHM

This sections describe the optimal algorithm we introduced in the main paper. Algorithm 1 shows the pseudocode of our algorithm. The optimal algorithm is a best-first search guided by our proposed heuristic and consists of the following steps:

1. Compute the exact quantities (Section 2.2) for every concept in the dataset (line 5).
2. For each concept, compute $dIoU_{max}$ and $dIoU_{min}$ for all possible paths starting from it, using the heuristic aggregated computation (line 6).
3. Initialize the frontier with all paths (line 7) and keep track of the greatest $dIoU_{min}$ (line 8). The frontier is sorted by the estimated $dIoU_{max}$.

Table 6: Avg. Time across 50 units (min) per hyperparameters on moderate settings.

Value	Our	MMESH
Explanation Len		
3	0.16	0.62
5	0.16	1.28
10	0.17	4.55
20	0.18	> 240
Beam Size		
5	0.16	0.62
10	0.17	1.24
20	0.18	> 240

4. Reduce the frontier by removing nodes whose estimated $dIoU_{max}$ is lower than the global maximum of the minimum estimates (line 10).
5. Iteratively pop nodes from the frontier, starting from those with the highest estimated $dIoU$ (lines 11-12).
 - (a) If the estimate is aggregated, refine it by computing the sample-based estimate and reinsert the node (lines 13-20). Otherwise, proceed to the next step.
 - (b) Apply logical equivalence rules when possible (lines 21-26), recompute the sample-based estimate for the equivalent expression, and reinsert the node if the estimate has changed. Otherwise, proceed to the next step. Currently, we only check distributive properties as logical equivalences. Note that even though the expressions are logically equivalent, their estimates may differ due to overestimation. The goal of this step is to select the form with the smallest overestimation among all possible equivalents.
 - (c) Check whether the same node has been recently explored by comparing its maximum IoU with the most recent one. If they match but the node has not yet been explored, add it to memory; if it has already been explored, skip it. Otherwise, clear the memory and initialize the most recent IoU with the node’s maximum IoU (lines 27–38).
 - (d) If the node is not a final one (i.e., can still be extended), expand its label by concatenating every possible concept with every allowed connective (line 49). For each new label, compute $dIoU_{max}$ and $dIoU_{min}$ via the heuristic aggregated computation and insert them into the frontier (lines 50-51). If a new maximum is found among $dIoU_{min}$ of the new paths, update the global maximum and reduce the frontier (lines 52-55).
 - (e) If the node is a final one, compute its exact IoU (line 42). During this step, the algorithm stores the intermediate quantities of the sub-labels composing the label (line 40). This information is then backpropagated to the frontier: nodes that share sub-labels with the evaluated node update their estimates using these exact quantities (line 41). If the IoU is greater than that of the best label found so far, update the best label with the current node (lines 43–45).
 - (f) Continue until the frontier is empty (line 11).

In the following, we provide the rationale behind several design choices made during the development of this algorithm.

Memory Mechanism The memory mechanism (lines 27–38) was introduced to account for logical equivalences not handled elsewhere in the algorithm and to prevent redundant computation. We considered several alternatives, but this solution proved to be the least computationally expensive. Without such a mechanism, some logically equivalent rules could be expanded multiple times (lines 49–55), leading to a significantly larger search space. Alternative options would be to check for logical equivalences or the presence of a node directly when adding nodes to the frontier. However, this would make exploration prohibitively expensive: verifying logical equivalence or membership would offset the efficiency gains of the heuristic. In addition, since the frontier is implemented as

864 a heap, checking whether a node is already present is slower than with a sorted list. Maintaining a
865 sorted frontier would require re-sorting on every insertion, which would be too costly and therefore
866 infeasible.

867
868 **Concept Quantities** A key design choice concerns the handling of concept quantities. We store
869 quantities only for atomic concepts and do not cache them when estimating (lines 6, 14, and 50) or
870 computing the $dIoU$ (line 40). Consequently, the intermediate quantities of labels are recomputed
871 multiple times during the search. This decision was made for two main reasons. First, storing ad-
872 ditional quantities increases access time, which is critical because the algorithm frequently accesses
873 these values; this overhead could easily exceed the time required to recompute them from scratch.
874 Second, caching more quantities increases memory usage, potentially limiting the ability to run
875 parallel processes. Given the substantial runtime required for the most complex datasets, preserv-
876 ing the ability to parallelize is essential. Moreover, avoiding extensive caching keeps the approach
877 lightweight in terms of memory and resource requirements. In conclusion, the marginal gains from
878 more precise estimations are outweighed by the overhead, especially given that the algorithm tends
879 to converge toward a breadth-first search (Section 3.3).

880 **Backpropagation** This process was introduced to reduce the number of visited states. While it has
881 no (or bad) impact on low and moderate-complexity settings, it significantly reduces runtime in high
882 complexity scenarios, especially when running on CPUs. Specifically, the frontier in the later stages
883 of the search often contains similar explanations that differ only in the last added concept (e.g., (A
884 AND B) OR C) and (A AND B) OR D). These nodes typically appear in the frontier because the left-
885 hand terms provide a rough overestimation relative to the current maximum. When a node is visited,
886 this overestimation is temporarily corrected, and the backpropagation step updates the estimates of
887 all remaining nodes. This correction helps the algorithm to avoid visiting unpromising states. On
888 average, 2000–3000 nodes per unit are updated via backpropagation, highlighting the importance of
889 this mechanism.

890 **Sample vs Aggregated Computation** The optimal algorithm leverages both sample-based and
891 aggregated computations for path and label estimations. This design is necessary due to the large
892 state space and the overestimation introduced by the aggregated computation. We explored several
893 alternatives during the development of this work, but this trade-off is the only one that ensures fea-
894 sibility in high complexity settings. As explained in the main text, sample computation requires $|\mathcal{D}|$
895 comparisons per calculation, whereas aggregated computation relies solely on the sum of concept
896 quantities computed at the first step of the algorithm, combining them in a single operation. These
897 two approaches represent a trade-off between precision and efficiency: the sample version is more
898 accurate but computationally intensive, while the aggregated version is faster but less precise. In
899 practice, using only aggregated computation for both node expansion and frontier exploration would
900 yield faster per-node computations but result in a frontier that is orders of magnitude larger than that
901 explored by our combined approach. Conversely, using only sample computation slightly reduces
902 the frontier (by roughly 50,000 nodes per unit in preliminary experiments) but incurs significantly
903 higher computation time per node, effectively nullifying the gains from the smaller frontier.

904 **MinIoU** As explained previously, the algorithm computes both the minimum and maximum IoU
905 for each label and path. This increases the time required per estimation, since the maximum and
906 minimum calculations are distinct and rarely share terms, effectively doubling the computation time
907 per node. An alternative design could omit the MinIoU computation and explore only nodes whose
908 maximum IoU exceeds the best visited so far. However, without the MinIoU, it becomes impossible
909 to dynamically reduce the frontier at runtime: the frontier can only grow, and the only way to
910 remove a single node is to fully explore it. This would result in a frontier orders of magnitude larger
911 than the one explored by the proposed design, especially in the early phase of the search, when
912 the MinIoU is updated multiple times and allows significant pruning. Future work could explore
913 adaptive strategies that decide dynamically whether to compute MinIoU based on the search space
914 or external information, potentially improving overall runtime.

915
916 **Code Optimization** In addition to the previously mentioned design choices, we implemented
917 several code optimizations to avoid unnecessary computation of quantities and to handle logical
equivalences. For instance, we enforce an order on concept indices when chaining two consecutive

918 applications of the same operator during node expansion. For example, for a concept with index 10,
 919 it can only be chained with concepts having an index greater than 10, since all other combinations
 920 will be captured when expanding nodes with lower indices. The same rule applies to consecutive
 921 operators of the same type. For example, if we have (3 OR 15), we can only chain concepts with
 922 indices greater than 15 when applying the OR operator again (e.g., (3 OR 15) OR 18), while we
 923 are free to choose any concept for other operators (AND or AND NOT). Other code optimizations
 924 involve avoiding the computations of paths when the intersection of the right or left sides is 0 and
 925 avoiding the sample computation related to equations involving $Bott_1$ in datasets where this vector
 926 is always 0 (see Section E)

927 C.1 FACTORS IMPACTING THE RUNTIME

928 As discussed in the main text, the primary factor influencing runtime is the degree of neuron align-
 929 ment: the more aligned a neuron is, the shorter the runtime, since the heuristic can prune the search
 930 space more effectively. Beyond alignment, other factors can also affect runtime, including the num-
 931 ber and type of concepts and the resulting explanation length.

932 Explanation length matters because the exploration strategy empirically resembles a breadth-first
 933 search, due to the coarse approximation of the maximum possible improvement along a path. As a
 934 result, it inherits the typical limitations of breadth-first exploration. The type of concept also plays
 935 a role: overly general concepts that are widely distributed across the dataset (e.g., colors) tend to
 936 overlap with many concepts and many images, forcing the algorithm to compute expensive approx-
 937 imations for all of them. Similarly, when a neuron responds to a general concept (e.g., bird) but
 938 admits many specializations (e.g., different colors), these specializations often have highly similar
 939 IoUs, preventing effective pruning and therefore increasing runtime.

940 The same dynamic arises in extreme cases where a concept covers nearly the entire input across
 941 most samples. Fortunately, such broad concepts (e.g., color, shape, or very general categories) are
 942 few in number and rarely useful in practice when they are so widely shared within the dataset.

943 D BEAM SEARCH ALGORITHM

944 In this section, we provide a high-level description of the beam search guided by our heuristic and
 945 present its pseudo-code in Algorithm 2. This algorithm replaces the MMESH heuristic (La Rosa
 946 et al., 2023) with the label quantity estimates introduced in Section 2.2, within a standard heuristic-
 947 guided beam search framework. The procedure begins by following the initial two steps of the
 948 optimal algorithm but keeps only the best b concepts to form the initial beam. It then proceeds
 949 iteratively for n steps: at each iteration, the algorithm expands the nodes in the current beam, sorts
 950 the resulting state space according to the heuristic, and evaluates candidate nodes by computing
 951 their IoU . The top b nodes are then selected to form the next beam. The process terminates when no
 952 candidate improves the current best IoU^{max} , or when no nodes remain to be expanded or visited.
 953 Differently from the optimal algorithm, this algorithm relies solely on sample computations and
 954 does not make use of, or estimate, the paths introduced in Section 2.3.2. Beyond the advantages
 955 discussed in Section 3.1, our heuristic does not rely on spatial information and can therefore be
 956 applied across multiple domains without requiring modifications to the code or formulation.

957 E ESTIMATING QUANTITIES

958 This section presents the rationale and derivation of all estimations computed by our heuristic. We
 959 divide the discussion into per-sample estimations and aggregated computations.

960 E.1 SAMPLE COMPUTATION

961 E.1.1 LABEL QUANTITIES

962 In the following, we derive the estimations of the label quantities introduced in the main text for
 963 the sample-based computation. For clarity, we group the equations by operator and discuss each

operator separately. Because the quantities for the unique elements are exact and are directly derived from the definition and Observation 1, here we focus on the derivation of the common quantities.

$$|I_{min}^C(L)_x| = \max(|I_{min}^C(L_{\leftarrow})_x|, |I^C(L_{\rightarrow})_x|) \quad (22)$$

$$|I_{max}^C(L)_x| = \min(|I_{max}^C(L_{\leftarrow})_x| + |I^C(L_{\rightarrow})_x|, |N_x^C|) \quad (23)$$

$$|E_{min}^C(L)_x| = \max(|E_{min}^C(L_{\leftarrow})_x|, |E^C(L_{\rightarrow})_x|) \quad (24)$$

$$|E_{max}^C(L)_x| = \min(|E_{max}^C(L_{\leftarrow})_x| + |E^C(L_{\rightarrow})_x|, |SE_x^C|) \quad (25)$$

Derivation Equations (22) to (25) for the OR operator: We can start by noting that Equation (22) corresponds to the case where one side is a subset of the other. In this case, the equation gives the maximum cardinality of the intersection between the left and right sides, since the minimum elements are already included in the maximum and the OR is 1-preserving (i.e., the number of ones cannot be lower than before the combination). Conversely, Equation (23) corresponds to the case where the sides are disjoint in their extras, adjusted by the $|N_x^C|$ quantity. Indeed, because the left side is an overestimation, it may happen that the sum exceeds the limits, requiring readjustment using $|N_x^C|$. The estimation of the maximum and minimum extras follows the same reasoning, with the only difference that the common space extra $|SE_x^C|$ is used to adjust the quantity of the maximum extras.

$$|I_{min}^C(L)_x| = \max(|I_{min}^C(L_{\leftarrow})_x| + |I^C(L_{\rightarrow})_x| - |N_x^C|, 0) \quad (26)$$

$$|I_{max}^C(L)_x| = \min(|I_{max}^C(L_{\leftarrow})_x|, |I^C(L_{\rightarrow})_x|) \quad (27)$$

$$|E_{min}^C(L)_x| = \max(|E_{min}^C(L_{\leftarrow})_x| + |E^C(L_{\rightarrow})_x| - |SE_x^C|, 0) \quad (28)$$

$$|E_{max}^C(L)_x| = \min(|E_{max}^C(L_{\leftarrow})_x|, |E^C(L_{\rightarrow})_x|) \quad (29)$$

Derivation Equations (26) to (29) for the AND operator: In this case, the AND operator is 0-preserving. Therefore, when computing the minimum, we consider the scenario where a guaranteed overlap (i.e., both sides equal to 1) occurs. This happens when the sum of the two sides exceeds the maximum available space, represented by $|N_x^C|$ and $|SE_x^C|$, respectively. In such cases, the minimum guaranteed overlap is given by the difference between the sum and the cardinality of the available space. In all other cases, the best estimation we can provide is simply 0. The computation of the maximum corresponds to the case of fully overlapping concepts. Since the operator is 0-preserving, the equation in this case selects the minimum of the two sides for each sample.

$$|I_{min}^C(L)_x| = \max(|I_{min}^C(L_{\leftarrow})_x| - |I^C(L_{\rightarrow})_x|, 0) \quad (30)$$

$$|I_{max}^C(L)_x| = \min(|I_{max}^C(L_{\leftarrow})_x|, |N_x^C| - |I^C(L_{\rightarrow})_x|) \quad (31)$$

$$|E_{min}^C(L)_x| = \max(|E_{min}^C(L_{\leftarrow})_x| - |E^C(L_{\rightarrow})_x|, 0) \quad (32)$$

$$|E_{max}^C(L)_x| = \min(|E_{max}^C(L_{\leftarrow})_x|, |SE_x^C| - |E^C(L_{\rightarrow})_x|) \quad (33)$$

Derivation Equations (30) to (33) for the AND NOT operator: This operator behaves like the AND operator but combines a negated concept, flipping all the bits in the corresponding vectors. The maximum estimations follow the same reasoning as for the AND operator, except that in this case $|N_x^C| - |I^C(L_{\rightarrow})_x|$ and $|SE_x^C| - |E^C(L_{\rightarrow})_x|$ represent the bits that were originally 0 and are now 1, corresponding to the common intersection and the common extras of the negated concept. The minimum estimation corresponds to the case where the right side fully overlaps with the left side. In this case, due to the negation, all overlapping bits flip to 0. Since the AND operator is 0-preserving, this results in the loss of all left side elements shared with the right side.

E.1.2 PATH QUANTITIES

In the following, we derive the estimations of the path quantities introduced in the main text. For clarity, we group the equations by operator and discuss each operator separately. Note that we assume, for all the paths, that at least 1 concept is added to the label, since the case where no

concepts are added is covered by the “final path” obtained by applying Definition 2.5 to the label quantities. In all the following equations, t denotes the difference between the maximum length and the length of the label L .

$$|I_{min}(L)_x| = \max(|I_{min}^C(L)_x| + |I_{min}^U(L)_x|, Bott_1(I^C)_x + Bott_1(I^U)_x) \quad (34)$$

$$|I_{max}(L)_x| = \min(|I_{max}^C(L)_x| + Top_t(I^C)_x, |N_x^C|) + \min(|I_{max}^U(L)_x| + Top_t(I^U)_x, |N_x^U|) \quad (35)$$

$$|Union_{min}(L)_x| = |N_x| + \max(|E_{min}^C(L)_x| + |E_{min}^U(L)_x|, Bott_1(E^C)_x + Bott_1(E^U)_x) \quad (36)$$

$$\begin{aligned} |Union_{max}(L)_x| = & |N_x| + \min(|E_{max}^C(L)_x| + Top_t(E^C)_x, |SE_x^C|) \\ & + \min(|E_{max}^U(L)_x| + Top_t(E^U)_x, |SE_x^U|) \end{aligned} \quad (37)$$

Derivation of Equations (34) to (37) for the OR operator: Equation (34) and Equation (36) follow the same derivation as Equations (22) and (24). In this case, the added concept is represented by $Bott_1(I^C)$, which corresponds to the minimum intersection of any concept in a given sample and reflects the case of fully overlapping concepts. In practice, for most datasets, this quantity is always equal to 0; thus, it reduces to $|I_{min}^C(L)_x| + |I_{min}^U(L)_x|$. The same reasoning applies to $|Union_{min}(L)_x|$ and the extras. Note also that the label quantities can be lower than the bottom values, since they represent combinations of concepts that may further reduce these quantities. Finally, the maximum quantities correspond to the case where the concepts included in L are disjoint from those cumulated in the Top vectors. In this situation, the quantities are simply the sum of neuron activations, maximum quantities, and the Top vectors, adjusted by the available space $|N_x^U|$ and $|SE_x^U|$, respectively.

$$|I_{min}(L)_x| = 0 \quad (38)$$

$$|I_{max}(L)_x| = \min(|I_{max}^C(L)_x|, Top_1(I^C)_x) \quad (39)$$

$$|Union_{min}(L)_x| = |N_x| \quad (40)$$

$$|Union_{max}(L)_x| = |N_x| + \min(|E_{max}^C(L)_x|, Top_1(E^C)_x) \quad (41)$$

Derivation Equations (38) to (41) for the AND operator: This operator is simpler to derive. Specifically, Equations (26) and (28) corresponds to the case where the label and all the concepts included in the Top vectors are disjoint, and thus all the quantities reduce to 0. Conversely, Equation (40) and Equation (41) represent the case where the label fully overlaps with the quantities stored in Top_1 . Therefore, for the AND operator, only the quantities from the smaller concept are preserved per sample. Note that Top_1 is the only applicable Top vector here, since higher indices sum the cardinality of multiple concepts, which would violate the operations defined by the AND operator.

$$|I_{min}(L)_x| = |I_{min}^U(L)_x| \quad (42)$$

$$|I_{max}(L)_x| = |I_{max}^U(L)_x| + \min(|I_{max}^C(L)_x|, |N_x^C| - Bott_1(I^C)_x) \quad (43)$$

$$|Union_{min}(L)_x| = |N_x| + |E_{min}^U(L)_x| \quad (44)$$

$$|Union_{max}(L)_x| = |N_x| + |E_{max}^C(L)_x| + \min(|E_{max}^U(L)_x|, |SE_x^C| - Bott_1(E^C)_x), \quad (45)$$

Derivation Equations (42) to (45) for the AND NOT operator: Equations (30) and (32) corresponds to the same case as in Equations (26) and (28). However, since the AND NOT operator preserves all unique elements (Observation 1), the minimum intersection corresponds to the minimum unique intersection of the label, and the minimum union includes the minimum unique extras. Conversely, Equations (44) and (45) represents the case where the label fully overlaps with the negated concept, represented by $|N_x^C| - Bott_1(I^C)_x$ and $|SE_x^C| - Bott_1(E^C)_x$, respectively. As previously noted, in practice, for most datasets, $Bott_1$ is always equal to 0. Thus, the equations simplify to $|I_{max}^U(L)_x| + |I_{max}^C(L)_x|$ for the intersection and $|N_x| + |E_{max}^C(L)_x| + |E_{max}^U(L)_x|$ for the union, since by definition $|I_{max}^C(L)_x| < |N_x^C|$ and $|E_{max}^U(L)_x| < |SE_x^C|$ due to the fact that $|N_x^C|$ represents the total neuron common space and $|SE_x^C|$ represents the total available space for common extras.

E.2 AGGREGATED COMPUTATION

This section describes the aggregated computation of the common quantities. To improve readability, we shorten the notation $\sum_{x \in \mathcal{D}}$ to simply \sum , since there are no ambiguities and the summation is used exclusively for this purpose.

E.2.1 LABEL QUANTITIES

In the case of disjoint concepts, we can use the same equation as in the sample-based case described in Section 2.3.1, since these are already aggregated. For the common quantities, the modification consists of pre-computing the aggregate value per label rather than per sample. This requires computing the dataset-wide sum only for atomic concepts, while all higher-arity labels can be derived through arithmetic operations over these sums. Therefore:

$$|I_{min}^C(L)_x| = \begin{cases} \min(\sum |I_{min}^C(L_{\leftarrow})_x|, \sum |I^C(L_{\rightarrow})_x|) & \text{if } \oplus = OR \\ \max(\sum |I_{min}^C(L_{\leftarrow})_x| + \sum |I^C(L_{\rightarrow})_x| - |N^C|, 0) & \text{if } \oplus = AND \\ \max(\sum |I_{min}^C(L_{\leftarrow})_x| - \sum |I^C(L_{\rightarrow})_x|, 0) & \text{if } \oplus = AND NOT \end{cases} \quad (46)$$

$$|I_{max}^C(L)_x| = \begin{cases} \min(\sum |I_{max}^C(L_{\leftarrow})_x| + \sum |I^C(L_{\rightarrow})_x|, |N^C|) & \text{if } \oplus = OR \\ \min(\sum |I_{max}^C(L_{\leftarrow})_x|, \sum |I^C(L_{\rightarrow})_x|) & \text{if } \oplus = AND \\ \min(\sum |I_{max}^C(L_{\leftarrow})_x|, |N^C| - \sum |I^C(L_{\rightarrow})_x|) & \text{if } \oplus = AND NOT \end{cases} \quad (47)$$

$$|E_{min}^C(L)_x| = \begin{cases} \max(\sum |E_{min}^C(L_{\leftarrow})_x|, \sum |E^C(L_{\rightarrow})_x|) & \text{if } \oplus = OR \\ \max(\sum |E_{min}^C(L_{\leftarrow})_x| + \sum |E^C(L_{\rightarrow})_x| - |SE^C|, 0) & \text{if } \oplus = AND \\ \max(\sum |E_{min}^C(L_{\leftarrow})_x| - \sum |E^C(L_{\rightarrow})_x|, 0) & \text{if } \oplus = AND NOT \end{cases} \quad (48)$$

$$|E_{max}^C(L)_x| = \begin{cases} \min(\sum |E_{max}^C(L_{\leftarrow})_x| + \sum |E^C(L_{\rightarrow})_x|, |SE^C|) & \text{if } \oplus = OR \\ \min(\sum |E_{max}^C(L_{\leftarrow})_x|, \sum |E^C(L_{\rightarrow})_x|) & \text{if } \oplus = AND \\ \min(\sum |E_{max}^C(L_{\leftarrow})_x|, |SE^C| - \sum |E^C(L_{\rightarrow})_x|) & \text{if } \oplus = AND NOT \end{cases} \quad (49)$$

The derivation of these quantities follows the same rationale as the sample computation, and the only difference is the larger overestimation produced by the aggregation of label-wise quantities.

E.2.2 PATH QUANTITIES

Similarly to the sample-based computation, computing the path $dIoU$ requires estimating the maximum possible improvement. In this case, however, the Top^A and $Bott^A$ vectors are computed concept-wise rather than per sample. Specifically, for each quantity, we sort the values of all individual concepts in the dataset and compute cumulative sums into the Top^A and $Bott^A$ vectors, starting from the largest and smallest values, respectively. Substituting these into the equations of

the sample-based computation, we obtain:

$$I_{min}(L) = \begin{cases} \max(\sum |I_{min}^C(L)_x| + \sum |I_{min}^U(L)_x|, \\ \quad Bott_1^A(I^C) + Bott_1^A(I^U)) & \text{OR Path} \\ 0 & \text{AND Path} \\ \sum |I_{min}^U(L)_x| & \text{AND NOT Path} \end{cases} \quad (50)$$

$$I_{max}(L) = \begin{cases} \min(\sum |I_{max}^C(L)_x| + Top_t^A(I^C), |N^C|) + \\ \quad \min(\sum |I_{max}^U(L)_x| + Top_t^A(I^U), |N^U|) & \text{OR Path} \\ \min(\sum |I_{max}^C(L)_x|, Top_1^A(I^C)) & \text{AND Path} \\ \sum |I_{max}^U(L)_x| + \min(\sum |I_{max}^C(L)_x|, \\ \quad |N^C| - Bott_1^A(I^C)) & \text{AND NOT Path} \end{cases} \quad (51)$$

$$|Union_{min}(L)| = |^1N| + \begin{cases} \max(\sum |E_{min}^C(L)_x| + \sum |E_{min}^U(L)_x|, \\ \quad Bott_1^A(E^C) + Bott_1^A(E^U)) & \text{OR Path} \\ \max(\sum |E_{min}^C(L)_x| + Bott_1^A(E^C) - |SE^C|, 0) & \text{AND Path} \\ \sum |E_{min}^C(L)_x| & \text{AND NOT Path} \end{cases} \quad (52)$$

$$|Union_{max}(L)| = |^1N| + \begin{cases} \min(\sum |E_{max}^C(L)_x| + Top_t^A(E^C), |SE^C|) \\ \quad + \min(\sum |E_{max}^U(L)_x| + Top_t^A(E^U), |SE^U|) & \text{OR Path} \\ \min(\sum |E_{max}^C(L)_x|, Top_1^A(E^C)) & \text{AND Path} \\ \sum |E_{max}^C(L)_x| + \min(\sum |E_{max}^U(L)_x|, |SE^C| - Bott_1^A(E^C)) & \text{AND NOT Path} \end{cases} \quad (53)$$

The derivation of these quantities follows the same rationale as the sample computation case, and the only difference is that these represent a larger overestimation. However, note in this case $|Union_{min}(L)|$ can be greater than 0, unlike in the sample computation where $\forall x \in \mathcal{D}$, $Bott_1(E^C)_x = 0$, except in the rare degenerate case (not observed in any of the datasets tested in this paper) where a single sample contains all concepts in the dataset.

E.3 PATH COMBINATIONS

As in the previous section, we shorten here the notation $\sum_{x \in \mathcal{D}}$ to simply \sum . As mentioned in the main text, to estimate values for paths involving multiple operators, we select the maximum and minimum of each quantity across the operators considered. For example, when both OR and AND can appear along a path, we compute:

$$dIoU_{max}(OR, AND)_x = \frac{\max(\sum |I_{max}^{OR}(L)_x|, \sum |I_{max}^{AND}(L)_x|)}{\min(\sum |Union_{min}^{OR}(L)_x|, \sum |Union_{min}^{AND}(L)_x|)} \quad (54)$$

and

$$dIoU_{min}(OR, AND)_x = \frac{\min(\sum |I_{min}^{OR}(L)_x|, \sum |I_{min}^{AND}(L)_x|)}{\max(\sum |Union_{max}^{OR}(L)_x|, \sum |Union_{max}^{AND}(L)_x|)} \quad (55)$$

These expressions simplify to:

$$dIoU_{max}(OR, AND)_x = \frac{\sum |I_{max}^{OR}(L)_x|}{\sum |Union_{min}^{AND}(L)_x|} \quad (56)$$

and

$$dIoU_{min}(OR, AND)_x = \frac{\sum |I_{min}^{AND}(L)_x|}{\sum |Union_{max}^{AND}(L)_x|} \quad (57)$$

We can further rewrite them as:

$$dIoU_{max}(OR, AND)_x = \frac{\min(|I_{max}^C(L)_x| + Top_t(I^C)_x, |N_x^C|) + \min(|I_{max}^U(L)_x| + Top_t(I^U)_x, |N_x^U|)}{|N_x|} \quad (58)$$

and

$$dIoU_{min} = 0 \quad (59)$$

Now, let us consider the case where we explicitly design this combined quantity. As before, we observe that $dIoU_{min} = 0$. For the numerator, however, we can refine the estimation: including an AND operator at any step will, by design, remove all the unique quantities of the current label. Thus:

$$dIoU_{max}(OR, AND) = \frac{\min(|I_{max}^C(L)_x| + Top_t(I^C)_x, |N_x^C|) + Top_t(I^U)_x}{|N_x| + (|E_{max}^C(L)_x| + |Bott_1(E^C)_x| - |SE_x^C|)} \quad (60)$$

However, since $|Bott_1(E^C)_x|$ is 0 in most datasets, the denominator reduces to $|N_x|$. Thus, the only practical difference lies in the numerator, where the $|I_{max}^U(L)_x|$ term is missing. However, in general, $Top_t(I^U)_x$ is much larger than $|I_{max}^U(L)_x|$, since it includes the maximum per sample across all concepts in the dataset. This makes the difference between $dIoU_{max}^{(OR, AND)}$ (from explicit design) and our proposed simplification relatively small. Similar observations can be made for all the other combinations of the operators considered in this paper.

From a practical perspective, and given that the optimal algorithm often converges toward breadth-first exploration (Section 3.3), the gain of the refined design is marginal. Conversely, our simplified approach, based on combining the values of exclusive paths of operators, scales more efficiently and facilitates future extensions. Indeed, new logic operators can be incorporated into the optimal algorithm simply by providing their estimates for the exclusive path.

F PROOF OF OPTIMALITY

This section builds upon the derivations introduced in section E to demonstrate that the heuristics employed are admissible, ensuring the algorithm is guaranteed to find the optimal solution and is complete. This section assumes the reader to be familiar with the terminology introduced in Section 2.1.

F.1 PROOF OF ADMISSIBILITY OF THE LABEL QUANTITIES HEURISTIC

Let L be the label obtained by concatenating the label L_{\leftarrow} with the concept L_{\rightarrow} . This heuristic aims at estimating the IoU score of L . By Equation (4), the IoU score is equivalent to the decomposed $dIoU$ score. To be admissible, the heuristic cannot underestimate the numerator or overestimate the denominator. Thus, it must satisfy the following constraints:

$$\begin{cases} |I_{max}^U(L)_x| + |I_{max}^C(L)_x| \geq |I^U(L)_x| + |I^C(L)_x| & (61) \\ 0 \leq |E_{min}^U(L)_x| + |E_{max}^C(L)_x| \leq |E^U(L)_x| + |E^C(L)_x| & (62) \end{cases}$$

for all $x \in \mathcal{D}$. Equation 61 ensures that the heuristic returns an optimistic estimate of the intersection (numerator), while Equation 62 ensures a pessimistic estimate of the denominator of eq. (3).

We observe that unique elements are always disjoint and cannot be shared by definition 2.1. Therefore, their value can be computed exactly by using Observation 1, and thus $|I_{max}^U(L)_x| = |I^U(L)_x|$ and $|E_{min}^U(L)_x| = |E^U(L)_x|$.

Therefore, the constraints reduce to:

$$\begin{cases} |I_{max}^C(L)_x| \geq |I^C(L)_x| & (63) \\ |E_{min}^C(L)_x| \leq |E^C(L)_x| & (64) \end{cases}$$

where the lower bound in the second inequality is omitted because cardinalities are non-negative.

These constraints are satisfied by design because of the derivations described in the previous section. Namely, if L_{\leftarrow} and L_{\rightarrow} are disjoint, then they share no common elements, and the algorithm can compute the exact values. Thus $|I_{max}^C(L)_x| = |I^C(L)_x|$ and $|E_{min}^C(L)_x| \leq |E^C(L)_x|$. Note that in this case, the maximum and the minimum of these quantities coincide.

For the AND NOT operator, both the unique and the common elements reduce to those of the left side L_{\leftarrow} . Therefore, $dIoU(L) = dIoU(L_{\leftarrow})$, and any path reachable from L is also reachable from L_{\leftarrow} since they are equivalent.

Conversely, if L_{\leftarrow} and L_{\rightarrow} overlap, we analyze each operator \oplus .

Case OR operator. For the intersection, recall that

$$|I_{max}^C(L)_x| = \min(|I_{max}^C(L_{\leftarrow})_x| + |I^C(L_{\rightarrow})_x|, |N_x^C|) \quad (65)$$

If the sets overlap, then

$$|I^C(L)_x| = |I^C(L_{\leftarrow})_x| + |I^C(L_{\rightarrow})_x| - |I^C(L_{\leftarrow})_x \cap I^C(L_{\rightarrow})_x| \quad (66)$$

Because they overlap, they share at least one element, and thus

$$|I^C(L_{\leftarrow})_x \cap I^C(L_{\rightarrow})_x| \geq 1. \quad (67)$$

This implies

$$|I^C(L)_x| \leq |I^C(L_{\leftarrow})_x| + |I^C(L_{\rightarrow})_x| \quad (68)$$

Furthermore, by the definition of N_x^C ,

$$|I^C(L)_x| \leq |N_x^C| \quad (69)$$

Thus $|I_{max}^C(L)_x| \geq |I^C(L)_x|$.

For the extras, recall that

$$|E_{min}^C(L)_x| = \max(|E_{min}^C(L_{\leftarrow})_x|, |E^C(L_{\rightarrow})_x|) \quad (70)$$

If the sets overlap, then

$$|E^C(L)_x| = |E^C(L_{\leftarrow})_x \cup E^C(L_{\rightarrow})_x| \quad (71)$$

Because the cardinality of a set cannot exceed the cardinality of its union with another set, Equation (64) is satisfied regardless of which element is selected by the max operator.

Case AND operator. In this case, for the intersection,

$$|I_{max}^C(L)_x| = \min(|I_{max}^C(L_{\leftarrow})_x|, |I^C(L_{\rightarrow})_x|) \quad (72)$$

The true value is

$$|I^C(L)_x| = |I^C(L_{\leftarrow})_x \cap I^C(L_{\rightarrow})_x| \quad (73)$$

since the AND operator is 0-preserving. Because an intersection cannot exceed either operand, then

$$|I_{max}^C(L)_x| \geq |I^C(L)_x| \quad (74)$$

Regarding the extras, recall that

$$|E_{min}^C(L)_x| = \max(|E_{min}^C(L_{\leftarrow})_x| + |E^C(L_{\rightarrow})_x| - |SE_x^C|, 0) \quad (75)$$

The zero case is proven by the non-negativity of the cardinality. Conversely, if the max operator selects the first factor, then we can prove it never overestimates the true value by induction by fixing L_{\leftarrow} to an atomic concept. In this case, we have access to the exact computation and thus

$$|E_{min}^C(L)_x| = \max(|E^C(L_{\leftarrow})_x| + |E^C(L_{\rightarrow})_x| - |SE_x^C|, 0) \quad (76)$$

If the sets overlap, then the true value is

$$|E^C(L)_x| = |E^C(L_{\leftarrow})_x \cap E^C(L_{\rightarrow})_x| \quad (77)$$

Recall that by the definition of SE^C , any sum between the cardinality of an extra disjoint set is lower than the cardinality of SE^C . Formally,

$$|SE_x^C| \geq |E^C(L_j)_x| + |E^C(L_k)_x|, \forall L_j, L_k \in \mathcal{L}^1 |E^C(L_j)_x \cap E^C(L_k)_x = \emptyset \quad (78)$$

1296 We rewrite $|E^C(L_{\leftarrow})_x|$ and $|E^C(L_{\rightarrow})_x|$:

$$1297 \quad |E^C(L_{\leftarrow})_x| = |E^C(L_{\leftarrow})_x \setminus E^C(L_{\rightarrow})_x| + |E^C(L_{\leftarrow})_x \cap E^C(L_{\rightarrow})_x| \quad (79)$$

$$1299 \quad |E^C(L_{\rightarrow})_x| = |E^C(L_{\rightarrow})_x \setminus E^C(L_{\leftarrow})_x| + |E^C(L_{\leftarrow})_x \cap E^C(L_{\rightarrow})_x| \quad (80)$$

1301 Since the sets $E^C(L_{\leftarrow})_x \setminus E^C(L_{\rightarrow})_x$, $E^C(L_{\rightarrow})_x \setminus E^C(L_{\leftarrow})_x$, and $E^C(L_{\leftarrow})_x \cap E^C(L_{\rightarrow})_x$ are
1302 pairwise disjoint, then:

$$1303 \quad |E^C(L_{\leftarrow})_x \setminus E^C(L_{\rightarrow})_x| + |E^C(L_{\rightarrow})_x \setminus E^C(L_{\leftarrow})_x| + |E^C(L_{\leftarrow})_x \cap E^C(L_{\rightarrow})_x| \leq |SE_x^C| \quad (81)$$

1305 Combining Equation (81) with Equation (76) and the fact that the operator selected the first argu-
1306 ment, we obtain

$$1307 \quad |E^C(L_{\leftarrow})_x| + |E^C(L_{\rightarrow})_x| - |SE_x^C| = |E^C(L_{\leftarrow})_x \setminus E^C(L_{\rightarrow})_x| + |E^C(L_{\leftarrow})_x \cap E^C(L_{\rightarrow})_x| \\ 1309 \quad + |E^C(L_{\rightarrow})_x \setminus E^C(L_{\leftarrow})_x| + |E^C(L_{\leftarrow})_x \cap E^C(L_{\rightarrow})_x| - |SE_x^C| \quad (82)$$

1311 which is smaller than or equal to $|E^C(L_{\leftarrow})_x \cap E^C(L_{\rightarrow})_x|$ due to Equation (81), thus satisfying
1312 Equation (64). Since $|E_{min}^C(L)_x|$ does not overestimate the true value for atomic concepts, and
1313 any successive step uses this underestimation to compute $|E^C(L_{\leftarrow})_x| + |E^C(L_{\rightarrow})_x| - |SE_x^C|$, then
1314 Equation (64) holds for any label $L \in \mathcal{L}^n$.

1315 **Case AND NOT operator.** In this case, for the intersection,

$$1317 \quad |I_{max}^C(L)_x| = \min(|I_{max}^C(L_{\leftarrow})_x|, |N_x^C| - |I^C(L_{\rightarrow})_x|) \quad (83)$$

1318 The true value is

$$1320 \quad |I^C(L)_x| = |I^C(L_{\leftarrow})_x \setminus I^C(L_{\rightarrow})_x| = |I^C(L_{\leftarrow})_x \cap (N_x^C \setminus I^C(L_{\rightarrow})_x)| \quad (84)$$

1322 Since AND NOT flips the right side and is 1-preserving on the left:

- 1323 • If the minimum selects $|I_{max}^C(L_{\leftarrow})_x|$, then because a set difference cannot exceed its left
1324 operand and I_{max}^C is an overestimation,

$$1326 \quad |I_{max}^C(L_{\leftarrow})_x| \geq |I^C(L_{\leftarrow})_x \setminus I^C(L_{\rightarrow})_x| \quad (85)$$

- 1328 • If the minimum selects $|N_x^C| - |I^C(L_{\rightarrow})_x|$, then this equals $|N_x^C \setminus I^C(L_{\rightarrow})_x|$, and because
1329 intersections cannot exceed either operand,

$$1330 \quad |I_{max}^C(L)_x| \geq |I^C(L)_x| \quad (86)$$

1332 Regarding the extras, recall that

$$1333 \quad |E_{min}^C(L)_x| = \max(|E_{min}^C(L_{\leftarrow})_x| - |E^C(L_{\rightarrow})_x|, 0) \quad (87)$$

1335 The zero case follows from the non-negativity of cardinality. Conversely, if the max operator selects
1336 the first quantity, then we can prove it never overestimates the true value by induction, by fixing
1337 L_{\leftarrow} to an atomic concept, similarly to the AND case. In this case, we have access to the exact
1338 computation and thus

$$1339 \quad |E_{min}^C(L)_x| = \max(|E^C(L_{\leftarrow})_x| - |E^C(L_{\rightarrow})_x|, 0) \quad (88)$$

1341 If the sets overlap, then the true value is

$$1342 \quad |E^C(L)_x| = |E^C(L_{\leftarrow})_x \setminus E^C(L_{\rightarrow})_x| = |E^C(L_{\leftarrow})_x| - |E^C(L_{\rightarrow})_x \cap E^C(L_{\leftarrow})_x| \quad (89)$$

1344 This is clearly greater than $|E^C(L_{\leftarrow})_x| - |E^C(L_{\rightarrow})_x|$ and thus eq. (64) is satisfied. Similarly to
1345 the AND operator case, since $|E_{min}^C(L)_x|$ does not overestimate the true value for atomic concepts,
1346 and any successive step uses this underestimation to compute $|E_{min}^C(L_{\leftarrow})_x| - |E^C(L_{\rightarrow})_x|$, then
1347 Equation (64) holds for any label $L \in \mathcal{L}^n$.

1349 Thus, for all operators considered in this paper, **the heuristic never underestimates the true com-
mon intersection, and admissibility is satisfied.**

Aggregated Computation. The admissibility of the aggregated version follows directly from the observation that the estimation produced by aggregating quantities per label rather than per sample is always larger than the one computed per sample, because it assumes completely disjoint concepts for the OR and AND NOT operators and fully overlapping concepts for the AND operator (i.e., the extreme cases).

F.2 REMAINING MINIMUM ESTIMATIONS

While not necessary for the proof of admissibility, in this section we prove the minimum estimations for the quantities not discussed in the previous paragraphs, namely $I_{min}^C(L)_x$ and $E_{max}^C(L)_x$. This proof will be useful to establish the admissibility of the Path Heuristic and the completeness of the algorithm. Both follow the same reasoning used for $E_{min}^C(L)_x$ and $I_{max}^C(L)_x$, respectively.

Namely, the minimum estimations for the intersection used by the algorithm are

$$|I_{min}^C(L)_x| = \begin{cases} \max(|I_{min}^C(L_{\leftarrow})_x|, |I^C(L_{\rightarrow})_x|) & \text{if } \oplus = OR \\ \max(|I_{min}^C(L_{\leftarrow})_x| + |I^C(L_{\rightarrow})_x| - |N_x^C|, 0) & \text{if } \oplus = AND \\ \max(|I_{min}^C(L_{\leftarrow})_x| - |I^C(L_{\rightarrow})_x|, 0) & \text{if } \oplus = AND NOT \end{cases} \quad (90)$$

The true values are

$$|I^C(L)_x| = \begin{cases} |I^C(L_{\leftarrow})_x \cup I^C(L_{\rightarrow})_x| & \text{if } \oplus = OR \\ |I^C(L_{\leftarrow})_x \cap I^C(L_{\rightarrow})_x| & \text{if } \oplus = AND \\ |I^C(L_{\leftarrow})_x \setminus I^C(L_{\rightarrow})_x| & \text{if } \oplus = AND NOT \end{cases} \quad (93)$$

Following the same reasoning used in Section F.1 to prove the underestimation of E_{min}^C , we obtain that $|I_{min}^C(L)_x| \leq |I^C(L)_x|$ holds for the OR operator, because the cardinality of a set cannot exceed that of its union with another set.

Similarly, for the other operators we replace E^C with I^C and SE^C with N^C when needed. For the AND operator, we expand

$$|I^C(L_{\leftarrow})_x| + |I^C(L_{\rightarrow})_x| = |I^C(L_{\leftarrow})_x \setminus I^C(L_{\rightarrow})_x| + |I^C(L_{\leftarrow})_x \cap I^C(L_{\rightarrow})_x| + |I^C(L_{\rightarrow})_x \setminus I^C(L_{\leftarrow})_x| + |I^C(L_{\leftarrow})_x \cap I^C(L_{\rightarrow})_x| \quad (96)$$

and using that the three disjoint parts satisfy

$$|I^C(L_{\leftarrow})_x \setminus I^C(L_{\rightarrow})_x| + |I^C(L_{\rightarrow})_x \setminus I^C(L_{\leftarrow})_x| + |I^C(L_{\leftarrow})_x \cap I^C(L_{\rightarrow})_x| \leq |N_x^C| \quad (97)$$

we obtain $|I_{min}^C(L)_x| \leq |I^C(L)_x|$.

For the AND NOT operator, we can rewrite the true value as

$$|I^C(L)_x| = |I^C(L_{\leftarrow})_x \setminus I^C(L_{\rightarrow})_x| = |I^C(L_{\leftarrow})_x| - |I^C(L_{\leftarrow})_x \cap I^C(L_{\rightarrow})_x| \quad (98)$$

which is greater than $|I^C(L_{\leftarrow})_x| - |I^C(L_{\rightarrow})_x|$ and thus $|I_{min}^C(L)_x| \leq |I^C(L)_x|$.

Analogously, for E^C the estimations and the true values are

$$|E_{max}^C(L)_x| = \begin{cases} \min(|E_{max}^C(L_{\leftarrow})_x| + |E^C(L_{\rightarrow})_x|, |SE_x^C|) & \text{if } \oplus = OR \\ \min(|E_{max}^C(L_{\leftarrow})_x|, |E^C(L_{\rightarrow})_x|) & \text{if } \oplus = AND \\ \min(|E_{max}^C(L_{\leftarrow})_x|, |SE_x^C| - |E^C(L_{\rightarrow})_x|) & \text{if } \oplus = AND NOT \end{cases} \quad (99)$$

and

$$|E^C(L)_x| = \begin{cases} |E^C(L_{\leftarrow})_x \cup E^C(L_{\rightarrow})_x| & \text{if } \oplus = OR \\ |E^C(L_{\leftarrow})_x \cap E^C(L_{\rightarrow})_x| & \text{if } \oplus = AND \\ |E^C(L_{\leftarrow})_x \setminus E^C(L_{\rightarrow})_x| & \text{if } \oplus = AND NOT \end{cases} \quad (102)$$

Using the same reasoning as in Section F.1 for I_{max}^C , we see that $|E_{max}^C(L)_x| \geq |E^C(L)_x|$ holds for the AND operator because the cardinality of an intersection cannot exceed that of either operand. For the OR operator, $|E_{max}^C(L)_x| \geq |E^C(L)_x|$ holds because $|E^C(L)_x| \leq |E^C(L_{\leftarrow})_x| + |E^C(L_{\rightarrow})_x|$ and $|E^C(L)_x| \leq |SE_x^C|$. Finally, for the AND NOT operator the inequality follows because a set difference cannot exceed its left operand when the first term is selected, and because an intersection cannot exceed either operand when the second term is selected.

1404 F.3 PROOF OF ADMISSIBILITY OF THE PATH HEURISTIC

1405 Let L_i be a generic label of length i . This heuristic aims at estimating the maximum achievable
 1406 IoU score by any label L_j in the state space $\mathcal{L}_{L_i}^n$, encompassing all labels of length up to n with
 1407 $i < n$ and $t = n - i$, such that the leftmost part of each label L_j corresponds to L_i (i.e., $L_j =$
 1408 $((L_i \oplus L_k) \dots) \oplus L_l$).

1410 Let L_* be the label in $\mathcal{L}_{L_i}^n$ associated with the highest IoU score. To be admissible, the heuristic
 1411 must satisfy the following constraints:

$$1412 \quad \left\{ \begin{array}{l} |I^U(L_*)_x| + |I^C(L_*)_x| \geq |I^U(L_*)_x| + |I^C(L_*)_x| \\ 0 \leq |E^U(L_*)_x| + |E^C(L_*)_x| \leq |E^U(L_*)_x| + |E^C(L_*)_x| \end{array} \right. \quad (105)$$

$$1414 \quad \left\{ \begin{array}{l} |I^U(L_*)_x| + |I^C(L_*)_x| \geq |I^U(L_*)_x| + |I^C(L_*)_x| \\ 0 \leq |E^U(L_*)_x| + |E^C(L_*)_x| \leq |E^U(L_*)_x| + |E^C(L_*)_x| \end{array} \right. \quad (106)$$

1416 for all $x \in \mathcal{D}$. In the constraints and in the following, we use the term "true value" to refer to the ex-
 1417 act computation and use the symbol $\widehat{QUANTITY}(L_*)_x$ to express the estimation of the maximum
 1418 possible value reachable from L_i for the considered quantity. Equation 105 ensures that the heuris-
 1419 tic returns an optimistic estimate of the intersection (numerator), while equation 106 guarantees a
 1420 pessimistic estimate of the denominator of eq. (3).

1422 These constraints are satisfied by design due to the derivations described in the previous section
 1423 for the Top and Bott vectors, which encode the maximum and minimum possible improvements,
 1424 respectively. To prove this, let us consider the exclusive paths for each operator considered in this
 1425 paper.

1426 **Case OR Operator.** In this case, the heuristic estimation for both the numerator and denominator
 1427 in eq. (3) is

$$1429 \quad \widehat{|I^U(L_*)_x|} + \widehat{|I^C(L_*)_x|} = \min(|I_{max}^C(L)_x| + Top_t(I^C)_x, |N_x^C|) \\ 1431 \quad + \min(|I_{max}^U(L)_x| + Top_t(I^U)_x, |N_x^U|) \quad (107)$$

1432 and

$$1433 \quad \widehat{|E^U(L_*)_x|} + \widehat{|E^C(L_*)_x|} = \max(|E_{min}^C(L)_x| + |E_{min}^U(L)_x|, Bott_1(E^C)_x + Bott_1(E^U)_x) \quad (108)$$

1436 where we drop the $|N|$ term because it is common to both expressions and does not affect the proof.

1438 Without loss of generality, let us consider the ideal case where there exists a sequence of concepts up
 1439 to length n , $\{L_j, \dots, L_n\}$, such that each of them fully overlaps with L_i on the extras, and each of
 1440 them is pairwise disjoint and disjoint from the elements of L_i that already overlap with the neuron
 1441 activations. Moreover, assume that these concepts have the largest possible number of elements
 1442 overlapping with the neuron activation. This construction represents the ideal scenario for the OR
 1443 operator, and no other concept in the dataset can further improve the IoU.

1444 In this ideal case, the true value of the denominator is

$$1445 \quad |E^U(L)_x| + |E^C(L)_x| \quad (109)$$

1446 because all concepts fully overlap with L on the extras.

1448 By comparing the true value with Equation (108), we obtain

$$1449 \quad \widehat{|E^U(L_*)_x|} + \widehat{|E^C(L_*)_x|} \leq |E^U(L_*)_x| + |E^C(L_*)_x| \quad (110)$$

1451 Indeed, if the max operator selects the first quantity, the inequality holds because $|E_{min}^C(L)_x|$ and
 1452 $|E_{min}^U(L)_x|$ are underestimations, as proven in Section F.1. If the max operator selects the second
 1453 quantity, the inequality follows from the definition of the $Bott_1(E^C)_x$ and $Bott_1(E^U)_x$ vectors
 1454 (Section 2.3.2). These vectors store, for each sample x , the minimum number of extras among all
 1455 concepts in the dataset (including 0). Thus, for any concept L_k in $\{L_j, \dots, L_n\}$, if L_k has the
 1456 smallest number of common extras for that sample, then $Bott_1(E^C)_x = |E^C(L_k)_x|$; otherwise,

$$1457 \quad Bott_1(E^C)_x < |E^C(L_k)_x| \quad (111)$$

1458 and thus

$$1459 \quad Bott_1(E^C)_x \leq |E^C(L_*)_x| \quad (112)$$

1460 The same argument applies to unique extras, and therefore Equation (106) is satisfied.

1461
1462 In the numerator case, because all the concepts are pairwise disjoint on the intersection and disjoint
1463 with respect to L_i , we can sum their quantities for ease of understanding as a single aggregated
1464 factor:

$$1465 \quad |I^C(L_*)_x| = |I^C(L_j)_x| + \dots + |I^C(L_n)_x| \quad (113)$$

$$1466 \quad |I^U(L_*)_x| = |I^U(L_j)_x| + \dots + |I^U(L_n)_x| \quad (114)$$

1467
1468 The true value of the numerator in this ideal case is given by

$$1469 \quad |I^C(L_i)_x| + |I^U(L_i)_x| + |I^C(L_*)_x| + |I^U(L_*)_x| \quad (115)$$

1471 By comparing the true value with Equation (107), we can note that when the minimum selects $|N_x^C|$
1472 and $|N_x^U|$, we have

$$1473 \quad \widehat{|I^C(L_*)_x|} + \widehat{|I^U(L_*)_x|} \geq |I^C(L_*)_x| + |I^U(L_*)_x| \quad (116)$$

1474 by definition, since $\{L_j, \dots, L_n\}$ are disjoint in the intersection with respect to L_i and their sum
1475 cannot exceed the available space for the common and unique intersections.

1476
1477 If the minimum selects the first terms, then as proven in Section F.1, $|I_{max}^C(L_i)_x|$ and $|I_{max}^U(L_i)_x|$
1478 are overestimations, and thus the constraint reduces to proving that

$$1479 \quad Top_t(I^C)_x \geq |I^C(L_*)_x| \quad (117)$$

$$1480 \quad Top_t(I^U)_x \geq |I^U(L_*)_x| \quad (118)$$

1481
1482 This is satisfied by construction of the Top_t vectors. These vectors store, for each sample, the sum
1483 of the t concepts in the dataset with the largest number of intersection elements for that sample.
1484 Note that t is also the length of L_* . Thus, for any sample $x \in \mathcal{D}$, if the t concepts with the largest
1485 intersections correspond to those in L_* , then $Top_t(I^C)_x = |I^C(L_*)_x|$, and otherwise $Top_t(I^C)_x >$
1486 $|I^C(L_*)_x|$. The same argument applies to the unique intersection, and therefore Equation (105) is
1487 satisfied.
1488

1489 **Case AND Operator.** In this case, the heuristic estimation for both the numerator and denominator
1490 in eq. (3) is

$$1491 \quad \widehat{|I^U(L_*)_x|} + \widehat{|I^C(L_*)_x|} = \min(|I_{max}^C(L)_x|, Top_1(I^C)_x) \quad (119)$$

1492 and

$$1493 \quad \widehat{|E^U(L_*)_x|} + \widehat{|E^C(L_*)_x|} = 0 \quad (120)$$

1494 where we drop the $|N|$ term because it is common to both expressions and does not affect the proof.
1495 We can note that Equation (106) is satisfied by the non-negativity of the cardinality.

1496
1497 To prove Equation (105), without loss of generality, let us consider the ideal case where there exists
1498 a sequence of concepts up to length n , $\{L_j, \dots, L_n\}$, such that each of them fully overlaps with L_i
1499 on the common intersection. This construction represents the ideal scenario for the AND operator
1500 since the AND operator cannot increase the intersection, but this construction allows preserving all
1501 the already intersecting common elements of L_i .

1502 The true value of the numerator in this ideal case is given by

$$1503 \quad |I^C(L_i)_x| \quad (121)$$

1504 since by definition the AND operator removes all the unique elements and, in this ideal case, all the
1505 common intersection of L_i are preserved.

1506
1507 If the minimum selects the first term, then as proven in Section F.1, $|I_{max}^C(L_i)_x|$ and $|I_{max}^U(L_i)_x|$
1508 are overestimations, and thus the constraint reduces to proving that

$$1509 \quad Top_1(I^C)_x \geq |I^C(L_i)_x| \quad (122)$$

$$1510 \quad Top_1(I^U)_x \geq |I^U(L_i)_x| \quad (123)$$

This is satisfied by construction of the Top_1 vectors and the fact in this ideal case the added concepts fully overlap. These vectors store, for each sample, the largest number of intersection elements for that sample among the concepts in \mathfrak{L}^1 . Thus, since any concept L_k in $\{L_j, \dots, L_n\}$ is assumed to fully overlap with L_i in the intersection, it either has the largest number of common intersection for that sample and thus $Top_1(I^C)_x = |I^C(L_k)_x|$, or it is not the largest and thus

$$Top_1(I^C)_x > |I^C(L_k)_x| \quad (124)$$

Therefore

$$Top_1(I^C)_x \geq |I^C(L_*)_x| \quad (125)$$

and Equation (106) is satisfied. In the general case, if any of the added concepts L_k does not fully overlap, then the true value becomes

$$|I^C(L_i)_x \cap I^C(L_k)_x| \quad (126)$$

which is satisfied by the fact that the cardinality of an intersection cannot exceed that of either operand when L_k is the concept associated with the largest common intersection in that sample and always in the other cases.

Case AND NOT Operator. In this case, the heuristic estimation for both the numerator and denominator in eq. (3) is

$$\widehat{|I^U(L_*)_x|} + \widehat{|I^C(L_*)_x|} = |I_{max}^U(L)_x| + \min(|I_{max}^C(L)_x|, |N_x^C| - Bott_1(I^C)_x) \quad (127)$$

and

$$\widehat{|E^U(L_*)_x|} + \widehat{|E^C(L_*)_x|} = |E_{min}^U(L)_x| \quad (128)$$

As done previously, let us consider the ideal case where there exists a sequence of concepts up to length n , $\{L_j, \dots, L_n\}$, such that each of them fully overlaps with L_i on the extras, and each of them is pairwise disjoint and disjoint from the elements of L_i that already overlap with the neuron activations. This construction represents the ideal scenario for the AND NOT operator, and no other concept in the dataset can further improve the IoU. We can note that the negation of these concepts corresponds to the ideal case of the AND operator, since the negation fully overlaps in the intersection and is disjoint in the extras. The difference here is that all the unique elements of L_i are preserved.

Thus, the true value of the numerator in this ideal case is

$$|I^C(L_i)_x| + |I^U(L_i)_x| \quad (129)$$

since the AND NOT operator preserves all unique elements and preserves all common intersections of L_i .

As proven in Section F.1, $|I_{max}^U(L_i)_x|$ is an underestimation, and thus the constraint reduces to proving that

$$\min(|I_{max}^C(L_i)_x|, |N_x^C| - Bott_1(I^C)_x) \geq |I^C(L_i)_x| \quad (130)$$

If the minimum selects the first term, then, as proven in Section F.1, $|I_{max}^C(L_i)_x|$ is an overestimation, and thus Equation (105) is satisfied. If the minimum selects the second term, we must prove that

$$|N_x^C| - Bott_1(I^C)_x \geq |I^C(L_i)_x| \quad (131)$$

Recall that $Bott_1$ stores, for each sample x , the minimum number of common intersections among all concepts in the dataset (including 0). Therefore, the quantity $|N_x^C| - Bott_1(I^C)_x$ represents the maximum number of intersecting elements obtainable from the negation of these concepts. In other words, it represents the maximum intersection space coverable by 0-entries for that sample.

Thus, for any concept L_k in $\{L_j, \dots, L_n\}$, if L_k has the smallest number of common intersections for that sample, then its negation corresponds to $|N_x^C| - Bott_1(I^C)_x$, and because it is disjoint with respect to L_i , we have

$$|N_x^C| - Bott_1(I^C)_x = |I^C(L_i)_x| \quad (132)$$

Conversely, if L_k does not have the smallest number of common intersections for that sample, then

$$Bott_1(I^C)_x \leq |I^C(L_i)_x| \quad (133)$$

1566 and therefore

$$1567 \quad |N_x^C| - Bott_1(I^C)_x > |I^C(L_i)_x| \quad (134)$$

1568 so Equation (105) is satisfied.

1570 In the general case, if any of the added concepts L_k is not fully disjoint, then the true value becomes

$$1571 \quad |I^C(L_i)_x \setminus I^C(L_k)_x| \quad (135)$$

1573 where L_k is the concept with the largest overlap with L_i . Since $I^C(L_k)_x$ removes at least
1574 $Bott_1(I^C)_x$ from $I^C(L_i)_x$, and $|I^C(L_i)_x| \leq |N_x^C|$ by definition, Equation (105) still holds.

1575 Regarding the denominator, since the concepts are disjoint in the extras, the true value is

$$1576 \quad |E^U(L_i)_x| \quad (136)$$

1578 because the AND NOT operator preserves all unique elements. As proven in Section F.1,
1579 $|E_{min}^U(L_i)_x|$ is an underestimation, and thus Equation (106) is satisfied.

1581 **Aggregated Computation.** The admissibility of the aggregated version follows the same proof as
1582 in the sample version. The only difference is that the *Top* and *Bott* vectors are computed per concept
1583 rather than per sample. In this case, they naturally represent overestimations and underestimations
1584 since they represent the ideal choice of the concepts to add for each operator at the dataset level
1585 (e.g., *Top* represents the case where the top t concepts in the dataset for the intersection are pairwise
1586 disjoint and disjoint with L_i over the full dataset).

1588 F.4 REMAINING MINIMUM ESTIMATIONS

1589 While not necessary for the proof of admissibility, in this section we prove that the minimum esti-
1590 mations for the path discussed in the previous paragraphs are underestimations. This result will be
1591 useful for guaranteeing the completeness of the algorithm.

1593 Namely, the minimum estimations for the intersection used by the algorithm are

$$1594 \quad \widehat{|I_{min}(L^*)_x|} \begin{cases} \max(|I_{min}^C(L)_x| + |I_{min}^U(L)_x|, & (137) \\ Bott_1(I^C)_x + Bott_1(I^U)_x) & \text{OR Path} & (138) \\ 0 & \text{AND Path} & (139) \\ |I_{min}^U(L)_x| & \text{AND NOT Path} & (140) \end{cases}$$

1600 We can note that the minimum estimations for the AND and AND NOT paths are underestimations
1601 by non-negativity of the cardinality in the AND path and by the property that the AND NOT operator
1602 preserves all unique elements.

1603 For the OR path, if the maximum operator selects the first term, then this represents an underesti-
1604 mation since $|I_{min}^C(L)_x|$ and $|I_{min}^U(L)_x|$ are underestimations of $|I^C(L)_x|$ and $|I^U(L)_x|$, as proven
1605 in Section F.2, and the OR operator cannot reduce the intersection. If the max operator selects the
1606 second term, then $\widehat{|I_{min}(L)_x|} \leq |I(L)_x|$ because OR is 1-preserving and any concept L_k chained
1607 to L_i through an OR operator must satisfy

$$1608 \quad |I^C(L_k)_x| + |I^U(L_k)_x| \geq Bott_1(I^C)_x + Bott_1(I^U)_x \quad (141)$$

1610 by definition of the $Bott_1$ vectors.

1612 Regarding the maximum denominator, recall that the algorithm uses the following estimate

$$1613 \quad \widehat{|Union_{max}(L^*)_x|} = |N_x| + \begin{cases} \min(|E_{max}^C(L)_x| + Top_t(E^C)_x, |SE_x^C|) & (142) \\ + \min(|E_{max}^U(L)_x| + Top_t(E^U)_x, & (143) \\ |SE_x^U|) & \text{OR Path} & (144) \\ \min(|E_{max}^C(L)_x|, Top_1(E^C)_x) & \text{AND Path} & (145) \\ |E_{max}^C(L)_x| + \min(|E_{max}^U(L)_x|, & (146) \\ |SE_x^C| - Bott_1(E^C)_x) & \text{AND NOT Path} & (147) \end{cases}$$

1619

where we can drop $|\mathcal{N}_x|$ since it is shared by both the true values and the estimations. These estimates are specular to the ones used for estimating the maximum possible intersection. Therefore, the result follows by applying the same reasoning used for $I_{max}^C(L)_x$ and replacing I^C with E^C , I^U with E^U , and N^C with SE^C .

F.5 OPTIMALITY OF THE ALGORITHM

Because the heuristics used to explore the state space are both admissible (Sections F.1 and F.3), because the quantities used to prune the frontier are guaranteed to be underestimations of the real IoU reachable from that node (Sections F.2 and F.4), and because the algorithm explores exhaustively all the nodes in the search space with an estimated alignment greater than the found solution, the algorithm is **complete and guaranteed to find the combination of concepts** $(L_1 \oplus L_2 \oplus \dots \oplus L_n)$ **among the concepts in \mathcal{L}^n that captures the highest possible alignment (IoU)**.

G EXAMPLES OF EXPLANATIONS DIFFERENCE

This section presents examples of the differences between explanations computed by beam search and those obtained with the optimal algorithm. To visualize the alignment, we select the samples in the dataset where the explanation holds and where the neuron is active within the considered activation range. We sort these samples by the size of the intersection between the explanation area and the activation area to ease interpretation, and then visualize the top 4 samples, highlighting the activation area in blue. The examples are not cherry-picked; rather, they correspond to the first nine differing explanations for the last convolutional layer of each explained model (ResNet18 (He et al., 2016), AlexNet (Krizhevsky et al., 2012), and DenseNet (Huang et al., 2017)). All the models have been pretrained on the Place365 dataset. Note that due to the selection procedure, samples used for visualization of different explanations or different methods for the same unit may differ, and this is expected. Alternative selection procedures, such as randomly extracting samples where the neuron is active, would produce harder-to-understand visualizations due to the superposition of neurons Elhage et al. (2022); O’Mahony et al. (2023); Dreyer et al. (2024) and variability in the size of activation. We also stress that the visualization is provided only as a reference and contextualization of the highlighted differences in explanations and not as a means of comparison. Despite the good results in the visualization, as discussed in the main text, the **optimality is not related to visualization properties or interpretability** but is related to the optimal (i.e., highest possible) solution found in a search problem and, specifically in this context, to the property of identifying the specific combination of concepts that captures (or expresses) the highest absolute spatial alignment with the neuron’s activations.

H OTHER NEURON AND CONCEPT-BASED EXPLANATIONS

While this paper focuses on methods targeting neuron spatial alignment, the literature offers a wide range of explanation approaches that use concepts to derive explanations Casper et al. (2023); Gilpin et al. (2018) or to decode other types of neurons Hesse et al. (2025); Srinivas et al. (2025); Bau et al. (2020); O’Mahony et al. (2025); Bykov et al. (2023) and layer behaviors Gao et al. (2025); Bricken et al. (2023). Our work is inspired by this growing interest in understanding neuron behaviors and using concepts to build explanations. However, the similarities end there: these approaches belong to distinct families of interpretability methods that differ fundamentally from compositional explanations across multiple dimensions, such as their goal (e.g., alignment vs. correlation), scope (e.g., neuron clusters vs. layers), assumptions (e.g., access to internals and availability of masks), representation (e.g., logical vs. statistical explanations), and phase (post-hoc vs. ante-hoc). Because these families are not able to measure the spatial alignment between activations and the locations of concepts, they have traditionally been regarded as complementary rather than competing categories with respect to compositional explanations. Consequently, they are not compared against them, and no established protocols exist for such comparisons. Given these differences, and because our work is theoretical in nature and specifically focuses on guaranteeing the optimality of compositional explanations without altering their established formulation, we leave the exploration of connections and complementarities with other interpretability paradigms for future research.

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727



Figure 2: Alignment detected in units of a ResNet18 model by both the optimal and beam-search methods. The blue regions indicate areas of neuron activation within the considered range.

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781



Figure 3: Alignment detected in units of a ResNet18 model by both the optimal and beam-search methods. The blue regions indicate areas of neuron activation within the considered range.

1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835

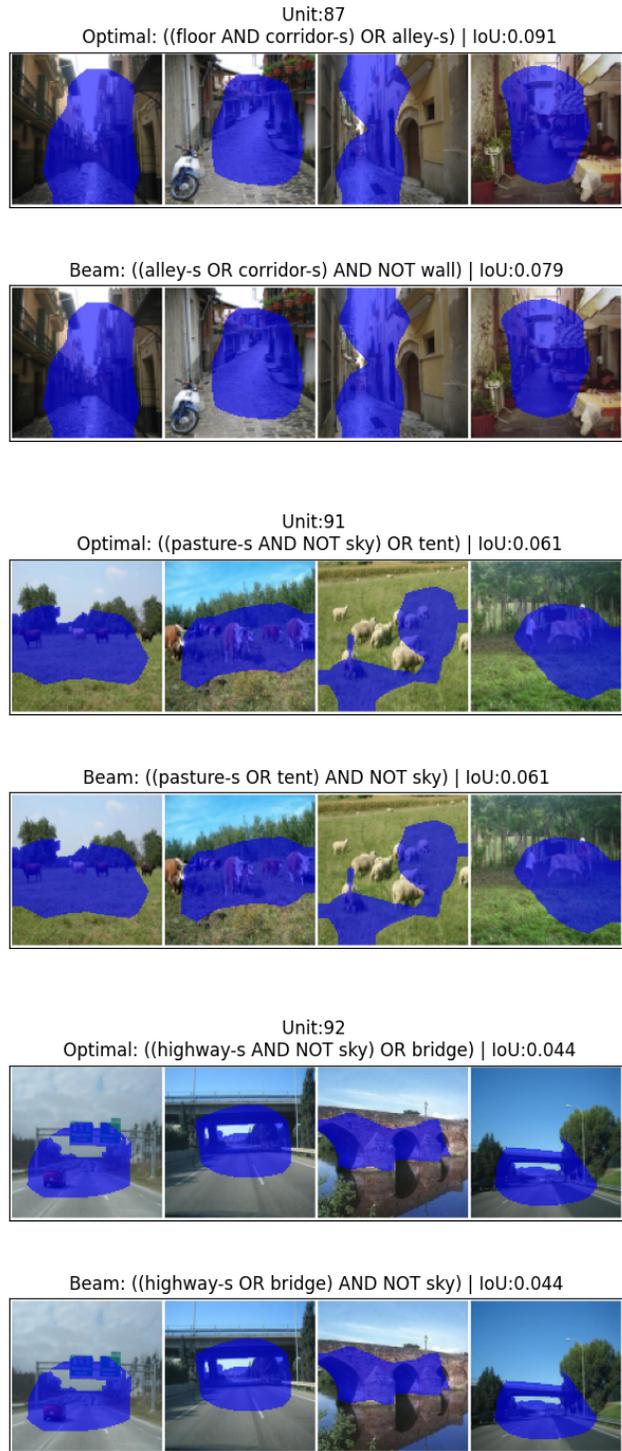


Figure 4: Alignment detected in units of a ResNet18 model by both the optimal and beam-search methods. The blue regions indicate areas of neuron activation within the considered range.

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

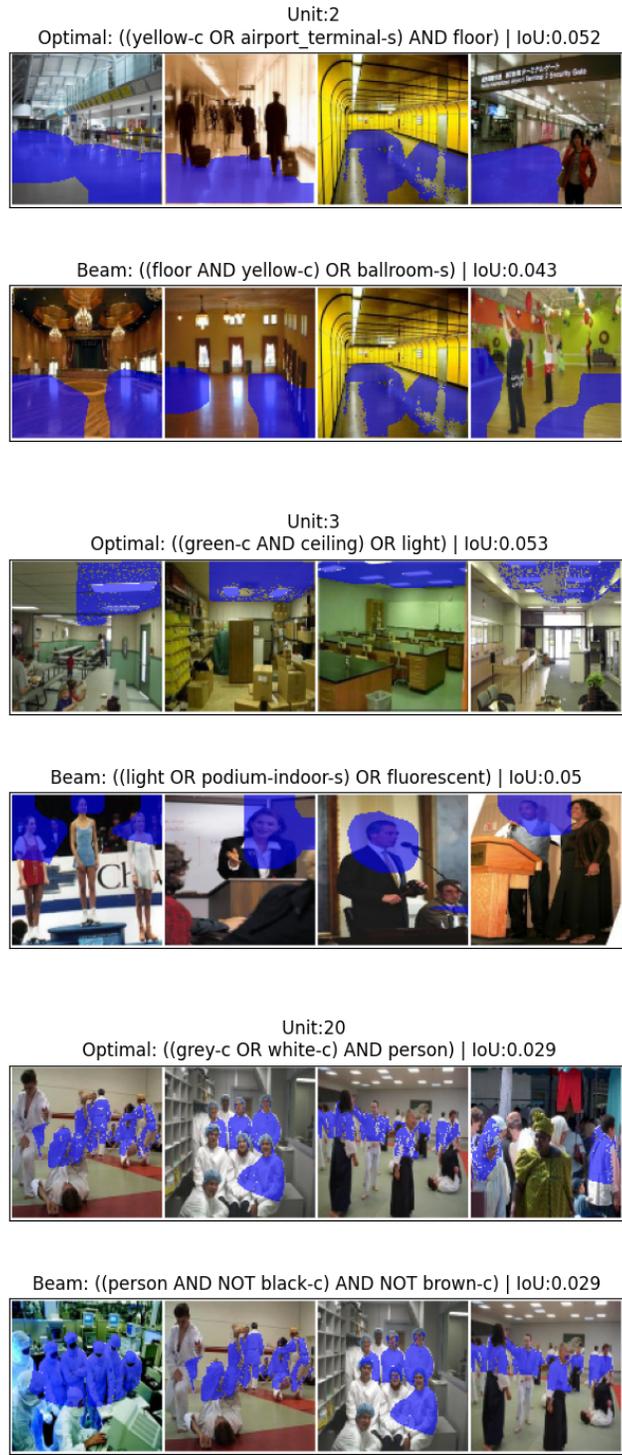


Figure 5: Alignment detected in units of an AlexNet model by both the optimal and beam-search methods. The blue regions indicate areas of neuron activation within the considered range.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

Unit:21
Optimal: ((waiting_room-s OR poolroom-home-s) AND floor) | IoU:0.031



Beam: ((road AND street-s) AND NOT white-c) | IoU:0.028



Unit:22
Optimal: ((bedroom-s OR living_room-s) AND ceiling) | IoU:0.047



Beam: ((ceiling AND living_room-s) AND NOT black-c) | IoU:0.041



Unit:26
Optimal: ((purple-c AND ceiling) OR pool table) | IoU:0.038



Beam: ((pool table OR ball_pit-s) OR day_care_center-s) | IoU:0.037



Figure 6: Alignment detected in units of an AlexNet model by both the optimal and beam-search methods. The blue regions indicate areas of neuron activation within the considered range.

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

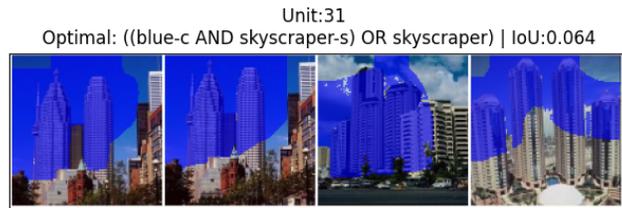


Figure 7: Alignment detected in units of an AlexNet model by both the optimal and beam-search methods. The blue regions indicate areas of neuron activation within the considered range.

1998
 1999
 2000
 2001
 2002
 2003
 2004
 2005
 2006
 2007
 2008
 2009
 2010
 2011
 2012
 2013
 2014
 2015
 2016
 2017
 2018
 2019
 2020
 2021
 2022
 2023
 2024
 2025
 2026
 2027
 2028
 2029
 2030
 2031
 2032
 2033
 2034
 2035
 2036
 2037
 2038
 2039
 2040
 2041
 2042
 2043
 2044
 2045
 2046
 2047
 2048
 2049
 2050
 2051

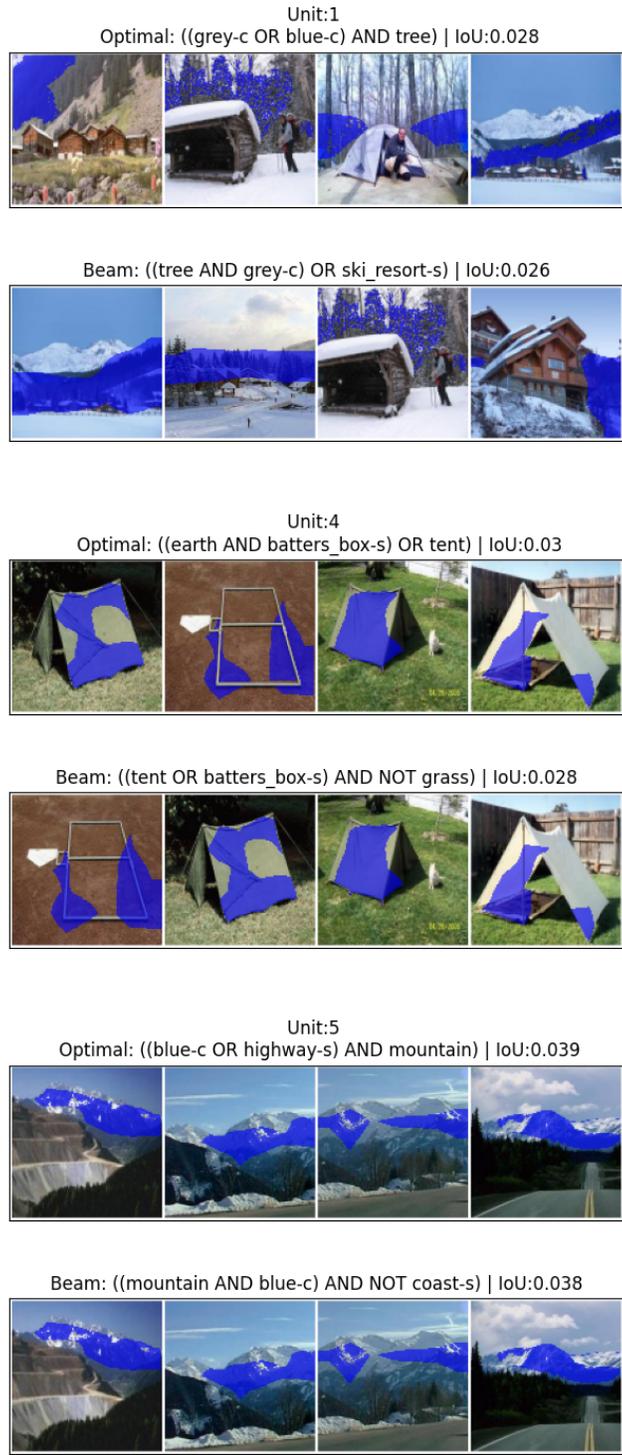


Figure 8: Alignment detected in units of a DenseNet161 model by both the optimal and beam-search methods. The blue regions indicate areas of neuron activation within the considered range.

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

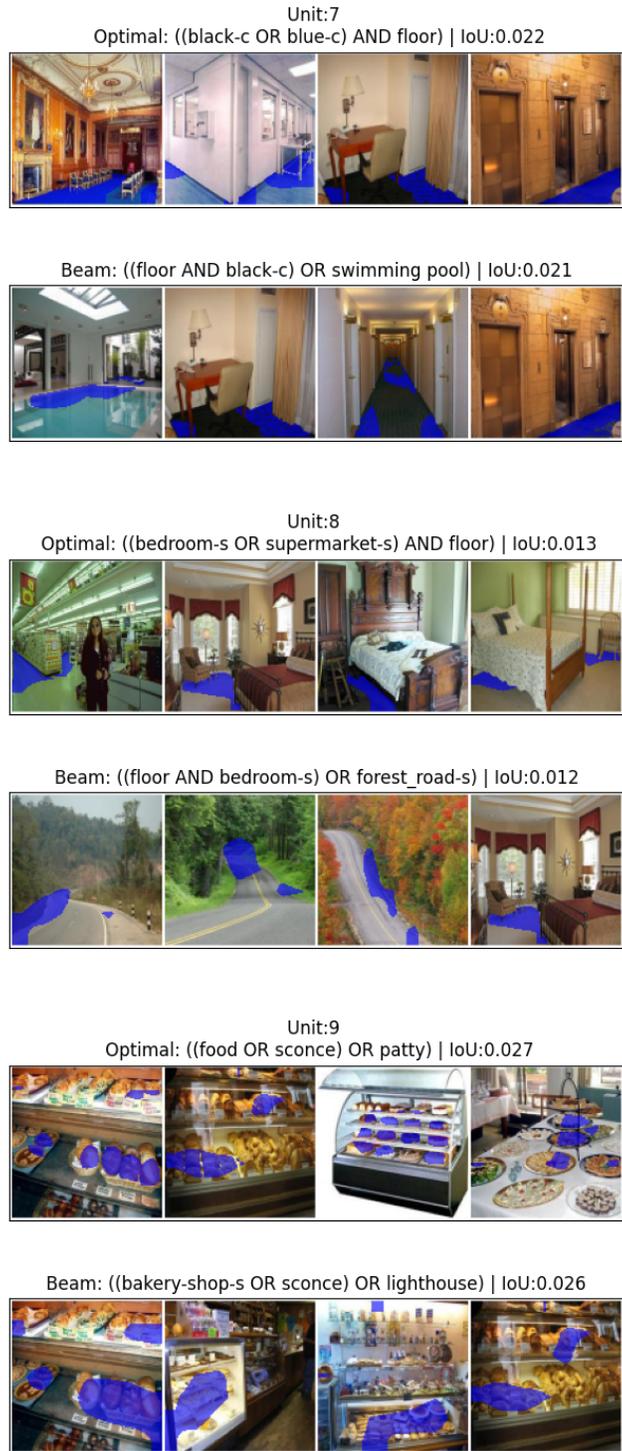


Figure 9: Alignment detected in units of a DenseNet161 model by both the optimal and beam-search methods. The blue regions indicate areas of neuron activation within the considered range.

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

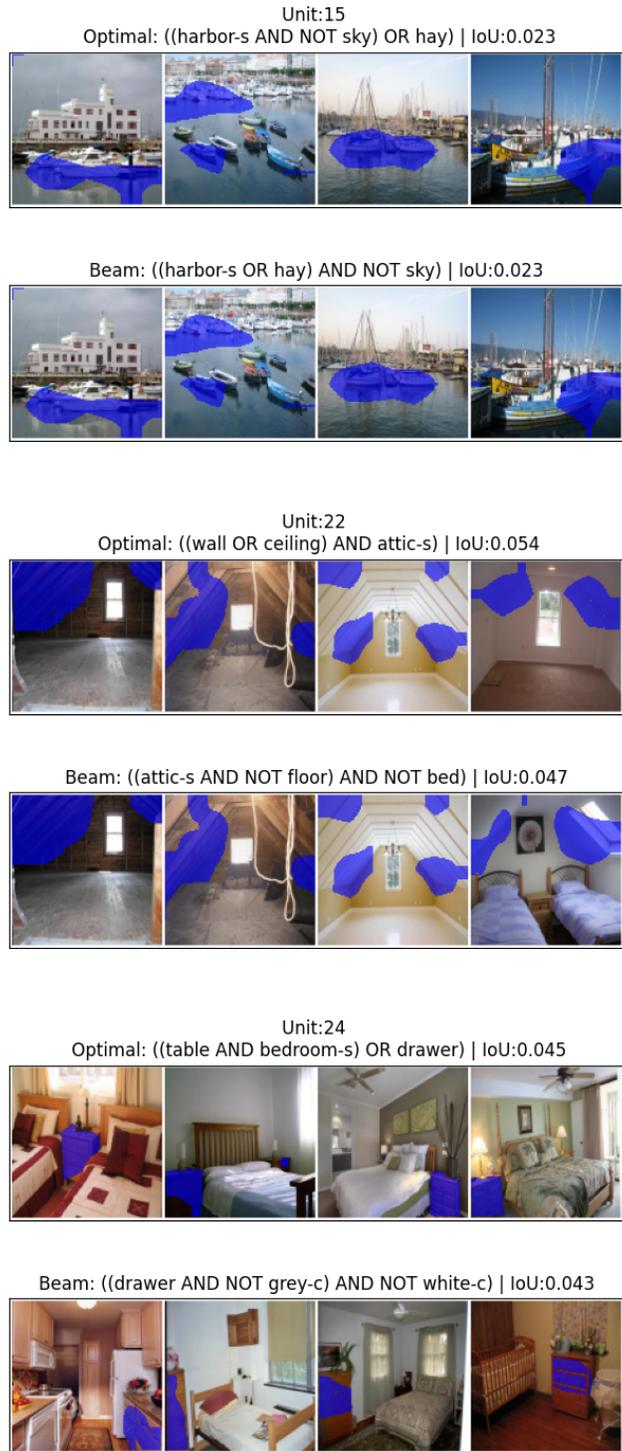


Figure 10: Alignment detected in units of a DenseNet161 model by both the optimal and beam-search methods. The blue regions indicate areas of neuron activation within the considered range.

2160 Conversely, our work is highly related to Network Dissection and Compositional Explanations. Net-
2161 work Dissection Bau et al. (2020) measures the spatial alignment between neurons and a set of in-
2162 dividual concepts, identifying the concept that maximizes alignment. Since it focuses on individual
2163 concepts, exhaustive search is feasible, and the resulting explanations are optimal within that con-
2164 text. Compositional Explanations extend this approach by searching for combinations of concepts
2165 that better express a neuron’s alignment. However, they are unable to explore the full state space,
2166 which sacrifices optimality. Our work closes this gap between Network Dissection and Composi-
2167 tional Explanations by providing a framework that guarantees optimal compositional explanations.
2168 Additionally, the development of our beam variant further improves runtime, narrowing the effi-
2169 ciency gap between Network Dissection and Compositional Explanations.

2170 2171 I LLM USAGE STATEMENT 2172

2173 In this work, large language models were used as a general-purpose assistive tool for polishing the
2174 text and grammar checking. Our workflow consisted of providing the LLM with a draft of some
2175 paragraphs and requesting it to highlight grammatical errors or suggest improvements. The output
2176 from the LLM was then used as a reference to improve the text, whenever the suggestions were
2177 deemed correct by the authors. The LLM did not contribute to research ideation or analysis at any
2178 stage.

2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

```

2214 Algorithm 1: Optimal Algorithm
2215
2216 Input:  $\mathcal{L}^1, N, \mathbf{M}$ , DisjointMatrix, length
2217 Output: BestLabel, BestIoU
2218 1 Frontier  $\leftarrow$  empty priority queue
2219 2 ConceptQuantities, Memory  $\leftarrow$  empty lists
2220 3 MinIoU, RecentIoU  $\leftarrow$  0
2221 4 for  $c_{k,i}$  in  $\mathcal{L}^1$  do
2222   5 | ConceptQuantities[ $c_{k,i}$ ]  $\leftarrow$  compute_quantities( $c_{k,i}, \mathbf{M}, N$ )
2223   6 | Paths  $\leftarrow$  estimate_aggregate_paths(ConceptQuantities[ $c_{k,i}$ ], length, MinIoU)
2224   7 | Frontier.add(Paths)
2225   8 | MinIoU  $\leftarrow$  update_min(Paths, MinIoU)
2226 9 end
2227 10 Frontier  $\leftarrow$  reduce_frontier(Frontier, MinIoU)
2228 11 while Frontier is not empty do
2229   12 | Node  $\leftarrow$  Frontier.pop()
2230   13 | if Node is aggregate_estimation then
2231     14 | UpdatedNode  $\leftarrow$  compute_sample_estimate(Node, MinIoU)
2232     15 | MinIoU  $\leftarrow$  update_min(UpdatedNode, MinIoU)
2233     16 | if UpdatedNode.max_iou > MinIoU then
2234       17 | | Frontier.add(UpdatedNode)
2235     18 | end
2236     19 | continue
2237   20 | end
2238   21 | UpdatedNode  $\leftarrow$  apply_logic_equivalences(Node)
2239   22 | if UpdatedNode.max_iou < Node.max_iou and UpdatedNode.max_iou > MinIoU then
2240     23 | | MinIoU  $\leftarrow$  update_min(UpdatedNode, MinIoU)
2241     24 | | Frontier.add(UpdatedNode)
2242     25 | | continue
2243   26 | end
2244   27 | if Node.iou == RecentIoU then
2245     28 | | if Node in Memory then
2246       29 | | | continue
2247     30 | | end
2248     31 | | else
2249       32 | | | Memory.add(Node)
2250     33 | | end
2251   34 | end
2252   35 | else
2253     36 | | Memory  $\leftarrow$  empty list
2254     37 | | RecentIoU  $\leftarrow$  Node.max_iou
2255   38 | end
2256   39 | if Node is final then
2257     40 | | TreeQuantities = compute_tree_quantities(Node)
2258     41 | | Frontier  $\leftarrow$  update_by_tree(Frontier, TreeQuantities)
2259     42 | | IoU = compute_iou(TreeQuantities.get(Node))
2260     43 | | if IoU > BestIoU then
2261       44 | | | BestIoU = IoU
2262       45 | | | BestLabel = Node.label
2263     46 | | end
2264   47 | | end
2265   48 | | else
2266     49 | | | AdditionalNodes  $\leftarrow$  expand(Node)
2267     50 | | | Paths  $\leftarrow$  estimate_aggregate_paths(AdditionalNodes, Quantities,
2268       51 | | | | length, MinIoU)
2269     52 | | | Frontier.add(Paths)
2270     53 | | | if min(Paths) > MinIoU then
2271       54 | | | | MinIoU  $\leftarrow$  min(Paths)
2272       55 | | | | Frontier  $\leftarrow$  reduce_frontier(Frontier, MinIoU)
2273     56 | | | end
2274   57 | | | end
2275   58 | | end
2276 end
2277 return BestLabel, BestIoU

```

2268
 2269
 2270
 2271
 2272
 2273
 2274
 2275
 2276
 2277
 2278
 2279
 2280
 2281
 2282
 2283
 2284
 2285
 2286
 2287
 2288
 2289
 2290
 2291
 2292
 2293
 2294
 2295
 2296
 2297
 2298
 2299
 2300
 2301
 2302
 2303
 2304
 2305
 2306
 2307
 2308
 2309
 2310
 2311
 2312
 2313
 2314
 2315
 2316
 2317
 2318
 2319
 2320
 2321

Algorithm 2: Our Informed Beam Search Algorithm

Input: \mathcal{L}^1 , N , M , DisjointMatrix, b , length

Output: BestLabel, BestIoU

```

1 Beam  $\leftarrow$  empty list
2 ConceptsQuantities  $\leftarrow$  empty list
3 for  $c_{k,i}$  in  $\mathcal{L}^1$  do
4   Quantities  $\leftarrow$  compute_quantities( $c_{k,i}$ ,  $N$ ,  $M$ )
5   ConceptsQuantities.append(Quantities)
6   IoU  $\leftarrow$  compute_dIoU(Quantities)
7   Beam.add(label =  $c_{k,i}$ , iou = IoU)
8 end
9 sort(Beam) # Sort by IoU
10 # Select the best b candidates
11 Beam  $\leftarrow$  Beam[:b]
12 MinIoU  $\leftarrow$  find_min(Beam)
13 for 2 to length do
14   SearchSpace  $\leftarrow$  expand_beam(Beam,  $\mathcal{L}^1$ )
15   Estimations  $\leftarrow$  estimate_labels_iou(SearchSpace, ConceptQuantities, ,
    DisjointMatrix)
16   sort(Estimations)
17   for L, EstIoU in Estimations do
18     if EstIoU < MinIoU then
19       # All the other labels cannot be added to the beam
20       break
21     end
22     iou  $\leftarrow$  compute_iou(L,  $N$ ,  $M$ )
23     Beam.add(label=L, iou=iou)
24   end
25   sort(Beam)
26   # Select the best b candidates
27   Beam  $\leftarrow$  Beam[:b]
28   # Compute and update info
29   MinIoU  $\leftarrow$  find_min(Beam)
30 end
31 BestLabel, BestIoU  $\leftarrow$  max(Beam)
32 return BestLabel, BestIoU

```
