Perception Matters: Enhancing Embodied AI with Uncertainty-Aware Semantic Segmentation

Sai Prasanna^{1*}, Daniel Honerkamp^{1*}, Kshitij Sirohi^{1*}, Tim Welschehold¹, Wolfram Burgard², and Abhinav Valada¹

¹Department of Computer Science, University of Freiburg, Germany ²Department of Engineering, University of Technology Nuremberg, Germany

Abstract. Embodied AI has made significant progress acting in unexplored environments. However, tasks such as object search have largely focused on efficient policy learning. In this work, we identify several gaps in current search methods: They largely focus on dated perception models, neglect temporal aggregation, and transfer from ground truth directly to noisy perception at test time, without accounting for the resulting overconfidence in the perceived state. We address the identified problems through calibrated perception probabilities and uncertainty across aggregation and found decisions, thereby adapting the models for sequential tasks. The resulting methods can be directly integrated with pretrained models across a wide family of existing search approaches at no additional training cost. We perform extensive evaluations of aggregation methods across both different semantic perception models and policies, confirming the importance of calibrated uncertainties in both the aggregation and found decisions. We make the code and trained models available at http://semantic-search.cs.uni-freiburg.de.

1 Introduction

Embodied AI has received a tremendous amount of attention in recent years, with new simulators enabling fast iteration in photorealistic, apartment-scale scenes based on real-world scans [15, 36]. This has been further advanced by institutionalized benchmarks and challenges in object search [1, 7, 40]. However, the main focus of this progress has primarily been on learning navigation and exploration policies. We analyze current object search works in Tab. 1, including, but not limited to, all top-performing approaches of the last three years' MultiOn, Habitat ObjectNav, and ImageNav challenges that released sufficient details. As most of these methods require a large number of episode steps, the overwhelming majority of approaches train with ground truth perception and then evaluate with pretrained out-of-the-box semantic perception models. This zero-shot transfer to a learned model can significantly reduce training costs. In return, we find it inducing a large gap between ground truth perception and semantic perception

^{*}These authors contributed equally.

Funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) - 417962828 and the Konrad Zuse School of Excellence in Learning and Intelligent Systems (ELIZA).

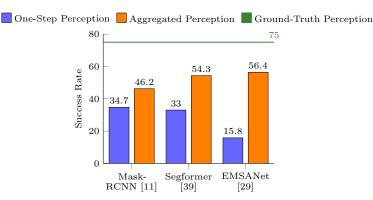


Fig. 1. Success rate of an RL agent [27] on the Habitat ObjectNav task with different semantic perception models. The gap from ground truth to learned perception models is often larger than the gap to an optimal policy. We propose uncertainty-based aggregation for sequential decision problems and find that this reduces the perception gap substantially. Ground Truth: ground truth semantic masks, One-Step: latest semantic prediction, Aggregated: best evaluated aggregation method of the model (cf. Sec. 4).

models, averaging over 25ppt, an error often as big as the gap to the optimal policy. Given their importance in the literature and the significance of the perception gap, this work focuses on the effective incorporation of such pretrained models.

We identify a number of unsolved problems in the literature: (i) Many approaches have noted false detections as major failure source [6,17,46]. However, due to their ease of use a vast share of the object search literature only evaluates comparatively old perception models. (ii) Object search differs from pure perception tasks as it requires decision making over a sequence of observations. This time dimension is commonly ignored, opting to simply use the latest or most likely prediction without any aggregation over time [5,42,46]. (iii) Zero-shot transfer from a perfect ground truth perception to an imperfect perception model signifies that policies are unaware that they act upon imperfect perception, which is worsened by the fact that the pre-trained perception models are generally overconfident in their predictions [30].

In order to quantify and address these problems, we first evaluate the impact of different semantic perception models and aggregation methods for sequential decision tasks. This differs from pure single-step perception evaluation based on the IoU (Intersection over Union) or precision [13]. In contrast, we measure the results over the full sequence of observations and actions, where early errors may impact or prevent later decisions. Fig. 1 shows the large perception gap to ground truth semantics. While newer models can reduce this gap, we find that temporal aggregation on the perception level is a key to closing the gap. To draw meaningful comparisons, we focus on one of the most used system structures, modular perception - mapping - policy pipelines, in one of the most explored tasks, ObjectNav. We introduce uncertainty-based aggregation and found decisions to address the identified problems. While previous methods develop complex, heuristic-based map aggregation strategies to cope with the overconfident and uncalibrated predicted probabilities [17,35], we incorporate calibrated perception

models with uncertainty estimation capabilities that can quantify this factor. In the second step, we evaluate these models with a learned search policy and across different semantic models. We find that our conclusions generalize to different search strategies and semantic perception models. The resulting methods directly integrate with existing approaches across a wide range of models without any additional training costs.

To summarize, this work makes the following main contributions:

- We identify and quantify gaps in the current integration of perception and aggregation methods in sequential embodied AI tasks.
- We incorporate calibrated perception probabilities and uncertainties in the agent's map aggregation and found decision making, thereby narrowing the gap to ground truth perception.
- We present extensive experimental evaluations and demonstrate that our approach can be easily integrated across a wide range of methods.
- We release the code at http://semantic-search.cs.uni-freiburg.de.

2 Related Work

Table 1. Comparison of existing object search approaches.

| Paper | Train | Test | Perception | Found | Suc | cess Ra | ate |
|------------------|------------------|----------------|--------------------------|---|-------------|-------------|-------------|
| | Perception | Perception | Aggregation | Decision | GT | Model | $_{ m Gap}$ |
| SemExp [5] | Mask-RCNN (PT) | Mask-RCNN (PT) | MaxPool | Distance | 73.1^{1} | 54.4^{1} | 18.7 |
| EEAux [42] | GT | RedNeT (FT) | _ | RL | 58.0^{2c} | 34.6^{2} | 23.4 |
| PONI [22] | GT | Mask-RCNN (FT) | Max-Pool | Distance | 86.5^{1} | 73.6^{1} | 12.9 |
| PONI [22] | GT | RedNeT (FT) | Max-Pool | Distance | 58.2^{3} | 31.8^{3} | 26.4 |
| Zhu et al. [47] | GT | Mask-RCNN (PT) | - | Distance | 76.8^{8} | 43.8^{8} | 33.0 |
| Stubborn [17] | - | RedNeT (FT) | Agg. Scores | Naive Bayes | 67.0^{2a} | 37.0^{2b} | 30.0 |
| Schmalstieg [27] | GT | GT | _ | GT | 85.7^{9} | _ | - |
| HIMOS [28] | GT | GT | _ | GT | 95.6^{9b} | _ | - |
| MOPA [25] | - | Detic (PT) | No Forgetting | RL (PT) | 81.0^{7} | 29.0^{7} | 52.0 |
| PIRLNav [24] | Demos | E2E | _ | RL | 61.9^{6} | _ | - |
| ESC [46] | - | GLIP (PT) | _ | Distance | 64.0^{6a} | 39.2^{6} | 24.8 |
| Dragon [6] | - | Mask-RCNN (PT) | Visual SLAM [†] | Distance | 84.2^{1} | 57.6^{1} | 26.6 |
| SkillFusion [35] | GT, part learned | SegFormer* | Binary w/ decay | $_{\substack{\text{Not-sure-action}}}^{\substack{\text{RL w}/\\ \text{Not-sure-action}}}$ | | 54.7^{5} | 10.0 |
| SkillTron [44] | GT, part learned | SegmATRon | Binary w/ decay | Distance | - | 59.0^{4} | - |
| Average | | | | | | | 25.8 |

GT: ground-truth information. PT: pretrained, FT: finetuned. Distance: if target is mapped, navigate to it, stop when close enough. No Forgetting: detected objects will never be deleted. Aggregated Scores: number of views, sum and max of class predictions. Demos: behavior cloning on expert demonstrations. Max-Pool: Channel-wise max-pooling of current and previous map. Part learned: GoalReacher trained with 20% noisy semantics. Not-sure action: the agent can decide to return to exploration if not confident to terminate the episode. Detection w/ decay: Accumulation of binary detections with decay coefficient. E2E: end-to-end learned, without explicit mapping. *no details provided if pretrained or trained from scratch. †no details on fusion provided.

Note that absolute success rates are not comparable across tasks. $^1\mathrm{Gibson}$ dataset, SemExp split. $^2\mathrm{MP3D}$, Habitat Challenge 2020 task, validation set, with 6 a or 15^b out of 21 target objects, c300 episode subset. $^3\mathrm{MP3D}$, Habitat Challenge 2021 task, validation set. $^4\mathrm{HM3D}$ v0.2, Habitat Challenge 2023 task, test set. $^5\mathrm{HM3D}$, 20 scenes of validation set. $^6\mathrm{HM3D}$, Habitat Challenge 2022, validation set, $^a\mathrm{Ground-truth}$ SR estimated from detection error rate. $^7\mathrm{HM3D}$, MultiOn 2.0. $^8\mathrm{MP3D}$ with 6 target categories and 10 test scenes (split unspecified). $^9\mathrm{iGibson}$ Challenge 2020, test set, $^b\mathrm{interactive}$ search.

Object Search has been tackled by a wide range of methods, including classical approaches such as frontier exploration [41], vision-based reinforcement learning [4], or auditory signals [43]. In recent years, map-based approaches that leverage semantic information have seen a large success [5,25,35,47] across public challenges [1]. On top of this map representation, both reinforcement learning

4

(RL) [5, 22, 28] and non-learned policies [17, 46] have been used successfully. This modularized network structure has found use in tasks beyond ObjectNav [3, 12]. We provide a comprehensive overview of recent approaches, including the top performing approaches of the last three year's ObjectNav [40], ImageNav [40] and MultiOn [38] challenges, in Tab. 1. We make a number of observations: (i) Due to the large size of perception models, learning-based approaches opt to use ground truth perception during training, to then replace this module with a pretrained perception model at test time. Occasionally, the perception model is first finetuned on the dataset [18]. (ii) Due to their ease of use, the majority of approaches rely on comparably old perception models such as Mask-RCNN [11] or Rednet [14]. Only a few works are able to stem the work to integrate recent state-of-the-art models [35,44,46]. (iii) While the policy networks often integrate previous information through RNNs, the time dimension of these sequential decision making processes is largely ignored for the perception level, and simply the last prediction is used. While separated works integrate aggregation methods [17,35], comprehensive comparisons within sequential decision tasks are missing. Outside this literature, methods employ techniques such as Bayesian update and probability averaging for map aggregation utilizing the non-calibrated probabilities from the perception semantic perception models [9, 19, 37, 48]. (iv) As a result, we find a large test-time gap between ground truth semantics and deployment with a learned model, averaging 25.8 ppt, with false detections often identified as the largest source of errors [6].

Uncertainty-Aware Perception: Uncertainty estimation techniques for perception can be classified into sampling-based and sampling-free methods. Sampling-based methods require either multiple forward passes [8] or multiple networks [16] to estimate the uncertainty. Therefore, they require significant time and memory and, thus, are not suitable for real-time applications. Hence, sampling-free methods have gained more attention in recent research. Sensoy et al. [30] proposed evidential learning for the classification task to learn parameters of a high-order distribution from which the classification uncertainty can be inferred. Other works have adapted evidential learning for object detection [21] and segmentation tasks with different modalities [32,33]. Sirohi et al. [31] further utilized the evidential output from [33] and proposed the methods for uncertainty-aware mapping. However, methods relying on evidential learning need to be trained from scratch. It is possible to directly utilize the softmax probabilities from the segmentation network in probabilistic log-odds softmax mapping [48]. However, the resulting map also creates over-confident uncertainty estimations [31]. Guo et al. [10] proposed temperature scaling to calibrate the softmax probabilities obtained from the network by scaling the predicted output logits on the validation set. This enables the use of any pre-trained network by tuning the scaling factor without the need for costly re-training. Hence, we utilize temperature scaling to obtain calibrated probabilities and uncertainties for semantic segmentation models.

3 Technical Approach

3.1 Problem Statement

In ObjectNav, the agent starts in an unexplored environment and has to find and navigate to an instance of a target category c, using only its RGB-D camera and localization. We follow the definition of the Habitat 2023 ObjectNav Challenge [40] in the HM3D dataset [23]. An episode is considered successful, if the agent issues a Stop action within 1 m of the target object and the object can be viewed by an oracle from the stopping location. The episode terminates unsuccessful if no found decision is made by 1,000 steps. We use the continuous action parameterization with some adaptations: (i) We flip the camera resolution from vertical to a more common landscape resolution of 640×480 and set the camera pitch to $-20 \deg$. (ii) We use only those target classes that are covered by all pretrained perception models, i.e. we omit the *plant* class. (iii) We adapt the step-size for continuous actions from $0.1 \,\mathrm{m\,s^{-1}}$ to $1.0 \,\mathrm{m\,s^{-1}}$, with a maximum linear velocity of $0.25 \,\mathrm{m\,s^{-1}}$. This ensures the agent can travel the same maximum distance in the allowed time budget as in the discrete action parameterization. (iv) We fix a bug that resulted in no collisions detection and low-level velocity integration to take multiple steps for larger velocities.

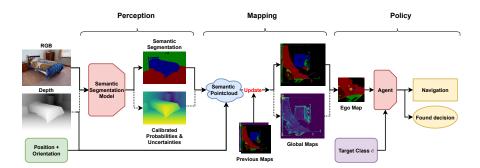


Fig. 2. Overview of modular object search pipelines. First, a semantic segmentation model classifies the current image. A mapping module then fuses this information into a semantic point cloud and integrates it into a global map. From this map, either an egocentric map is extracted for RL agents or the full map is used by a planner. An agent then determines navigation and found decision for a given target class c. We develop general methods to incorporate calibrated uncertainties in this system for temporal aggregation of the semantic perception and consistent found decisions.

3.2 Model Structure

In this work, we focus on modular approaches that have seen widespread use and achieved state-of-the-art results [5,6,17,22,25,27,28,35,44,46,47], depicted in Fig. 2. These approaches first obtain semantic information for the current camera image. They then use a mapping module to obtain an explicit representation, by projecting this data into a top-down map and integrating it into a global

map. Finally, an agent uses this representation to make navigation and found decisions - commonly via an ego-centric map for reinforcement learning agents, or using the full global map for planners. The modularity of this pipeline enables the use of different segmentation models with different policies. We incorporate uncertainty and calibrated probabilities into these systems and leverage them for efficient temporal aggregation and found decisions.

3.3 Uncertainty-Aware Perception

We assume access to a pretrained semantic segmentation network which takes RGB and/or depth images as input and predicts a semantic segmentation output as a logits vector $l = [l^1, ..., l^C]$ over the pixels, where C is the number of classes. However, these semantic segmentation models are commonly trained with a cross-entropy loss [14,29,39] which includes employing a softmax operation on the logits, which inflates the predicted probability of the class. Thus, we employ the temperature scaling technique as proposed by [10], where they tune a scaling factor on the validation set to scale the logits for better probability estimates. However, as the networks are pre-trained on different datasets, we tune the scaling factor t on a labeled set from the HM3D dataset [23]. Please note that this image set is independent of any policy training. Finally, we utilize the scaled logit vector $l_s = \left[\frac{l}{t}^1, ..., \frac{l}{t}^C\right]$ to obtain the probability vector $p_i = \operatorname{softmax}(l_s)$. Where, $p_i^{pred} = \left[p_1^1, ..., p_i^C\right]$ consists of per-class classification probability for the pixel i in the image. Finally, we employ normalized entropy to calculate the corresponding per-pixel classification uncertainty $u_i \in [0,1]$, defined by $u_i = \frac{\sum_{c=1}^C p_i^c \log(p_i^c)}{\log(C)}$.

3.4 Perception Uncertainty Weighted Map Aggregation

We use a Birds-Eye-View (BEV) grid map representation with a grid cell size of 3×3 cm². We project the perception prediction to the respective cell based on the depth image and the ground truth pose of the robot, using the top most voxel. We utilize the calibrated probability vector p^{pred} , together with the perception uncertainty u for the mapping. For every grid cell k in the map and N measurements, we calculate a weighted average of the probability vector, weighted by the inverse of the perception uncertainty, to obtain the aggregated probability vector $p_k = \frac{1}{U} \sum_{n=1}^N \frac{1}{u_{n,k}} p^{pred,k}$ of size $C \times 1$, where $U = \sum_{n=1}^N \frac{1}{u_{n,k}}$. Similar to the perception, we calculate the mapping uncertainty u_k^{map} for every grid cell as the normalized entropy obtained from the probability vector p_k . Thus, we maintain a vector m_k of size C+3 for every cell k of the map, where $m_k = [p_k, height, occupancy, u_k^{map}]$. The height keeps track of the maximum encountered height within a cell, and occupancy is a binary value, which we set to one if the height > 0.1 m, and zero otherwise.

3.5 Map Uncertainty-based Found Decision

We issue a found decision when the target object follows two criteria, (i) the target object is within 1 m of the robot, and (ii) the map uncertainty, u_k^{map} , of the map cell occupied by the target object is less than the threshold $\xi = 0.4$. We

set the value of the threshold empirically based on hyperparameter optimization on a training set (cf. Sec. 4). The idea is that perception is prone to make errors at farther distances. Hence, the distance constraint ensures the robot is in the vicinity of the target object when the decision is made. The map uncertainty is crucial to filter false positives caused by varying perception predictions. If the perception shows variation in the predictions, the uncertainty will be high. Thus, we only mark an object as found when perception provides multiple observations with low uncertainty for the target object.

3.6 Policies

While the perception literature has developed reliable metrics for single image evaluation, we are interested in the performance of perception models in sequential decision making problems. As such, the performance of perception and policy are tightly interwoven as the perceived state will impact the policy's next actions which in turn will alter the next state. As a result, some errors may have a larger impact than others, if they lead the agent astray early on or terminate the episode too soon. For this reason, we evaluate two settings:

First, we propose the use of a state-independent ground truth shortest path policy to evaluate all methods on the same sequence of observations to isolate the impact of the perception and aggregation components. We leverage Habitat's shortest path implementation [36] in the ground truth navigation mesh with the target object as the goal. For comparability, we collect metrics over the full trajectory until the goal, even if a false found decision would terminate the episode early.

In the second step, we evaluate how these results transfer to different policies and impact the overall performance within the decision making loop. For this, we implement a recent reinforcement learning based object search policy [27] that follows the modular model structure shown in Fig. 2. Given the target class, the agent transforms the full semantic map into a local and global ego-centric map and maps the task-relevant objects to a target-object color and non-relevant objects to an occupied color. Its policy network then predicts both the most likely direction toward the target object and navigation commands. The agent is trained with PPO with the ground truth semantic perception and then deployed with the learned perception model and aggregation strategy.

4 Experimental Evaluations

We evaluate our approach across a wide range of semantic perception models and multiple policies on the validation split of the Habitat 2023 ObjectNav Challenge [40] in the HM3D dataset [23].

Metrics: We compute the following evaluation metrics:

Success Rate (SR): The share of successful episodes in which the agent found the target and correctly raised the found decision within the time limit.

Found/False Positive Rate (FPR): Share of episodes with incorrect found decision.

Found/False Negative Rate (FNR): Share of episodes in which the agent failed to raise a found decision.

Detection/False Positives (#FP): Number of times the target object was incorrectly mapped and shown to the agent per episode. Note that multiple false detections can occur over a single episode. To count the number of object detections (as opposed to pixels), we dilate all target class predictions outside the ground truth bounding box on the map and count each connected component as a false detection.

Detection/False Negatives (#FN): Number of times per episode a target object existed in the map created with ground truth semantic camera, but no target class was shown to the agent within the ground truth bounding box of that object.

Success weighted by Path Length (SPL): Success weighted inverse normalized length of the agent's path.

Baselines: We compare against a large range of baselines for aggregation and found decision:

Ground Truth: Build the map with the ground truth semantic camera images from the simulator. Always trigger found decision if close enough to mapped target. Latest: Map the class with the highest predicted probability at each step, overwriting any previous values. Found decision is made if the agent is within the success radius $(1\,\mathrm{m})$ of the grid cell containing the goal object.

Hits/Views: Based on the probabilistic counting method [37] for aggregating multiple observations for binary classifications and views: Maintain additional map channels to track the number of hits (detections) of the goal class for each cell and how often this cell was in close view from a distance below d_{view} . If the agent is within the distance threshold of 1 m of a cell with at least v views and the hits/view ratio is above θ , raise the found decision. Otherwise, classify as false detection and stop mapping the target in that cell.

Skillfusion [35]: Erode the goal object in the local map with a kernel of 4 cm \times 4 cm to remove outliers. Then maintain a grid map with continuous values representing the existence of goals. This value is incremented by one if a goal object gets projected to the cell. Otherwise, the grid cell value is multiplied by a decay factor α . The goal is mapped for the agent if the cell value is greater than a threshold T. If the agent is within found distance of such a grid cell, raise the found decision. Stubborn [17]: Maintain additional map channels with total views, cumulative confidence, maximum confidence and maximum non-target confidence. These features are given to a Naive Bayes Classifier that outputs a binary found decision. As the trained classifier is not released, we train it on features collected from 64 episodes collected with the shortest path policy in unique training scenes.

Latest Filtered: Same as Latest, but only map the target category if the map uncertainty for the mapped target object is below a threshold ρ . Otherwise, map as occupied class.

Log Odds: Bayesian updating of the grid cells with a multi-class log odds vector for each cell [48] but using our calibrated probabilities. The agent is shown the most likely class at each step. If a goal object is within found distance and the uncertainty of the posterior is less than a threshold ξ , mark the object as found.

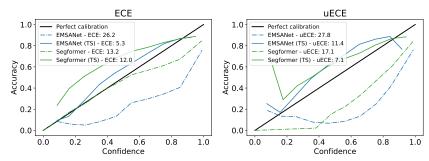


Fig. 3. Expected Calibration Error (left) and Uncertainty Expected Calibration Error (right) of the different semantic perception models on the validation set.

We set the values of the parameters for all methods via hyperparameter search with a tree-structured parzen estimator [2]. For each method, we set a budget of 20 trials and evaluate the value of each configuration as the average success rate over 30 episodes of the training scenes collected with the shortest path policy, and found the optimization process to converge to stable parameter values.

Perception Models: To evaluate the generalizability of the aggregation methods, we evaluate all approaches across different semantic models.

Mask-RCNN [11]: The most used perception model in the ObjectNav literature (cf. Tab. 1), representative of simple out-of-the-box models. Pretrained on MS-COCO. As Mask-RCNN is an instance segmentation model, it does not provide class probability for all pixels, hence not all aggregation methods are applicable. Segformer [39]: Transformer-based base model used by most recent ObjectNav challenge winners [35]. It is pretrained on ADE20k [45].

EMSANet [29]: A recent panoptic segmentation model which takes RGB and depth images as input. We only utilize semantic segmentation prediction, i.e., we only reason about the class of individual pixels and do not use the instance segmentation results of EMSANet. The model is pre-trained on the SUN RGB-D [34] and Hypersim [26] datasets.

4.1 Evaluation of Uncertainty Estimation

We evaluate the probability and uncertainty estimation quality of the perception models through the Expected Calibration Error (ECE) [20] and Uncertainty Expected Calibration Error (uECE) [33]. The calibration of a network defines how well the predicted confidence matches the actual accuracy of the prediction. ECE quantifies the error between the maximum class probability and the accuracy, whereas uECE measures the error between the network confidence and the accuracy, where confidence is defined as 1 – uncertainty. As the map aggregation as defined in the (Sec. 3.4) uses both uncertainty and probability, it is desirable to have both ECE and uECE to be low for better performance.

In Fig. 3, we show the calibration plots for vanilla and temperature-scaled versions of Segformer and EMSANet. A perfect calibration corresponds to the solid black line in the plot. As we can see, both ECE and uECE decrease after applying temperature scaling for both networks. While the ECE for EMSANet with

temperature scaling (TS) decreased from 26 percent to 5 percent, Segformer's only decreased by one percentage point. Thus the probabilities become much better calibrated with temperature scaling for EMSANet than for Segformer. Hence, we expect better average probability aggregation for EMSANet in comparison to Segformer. However, the uncertainty calibration of both EMSANet and Segformer shows a significant improvement in uECE of 17 and 10 percent respectively. Thus, uncertainty weighting should have an identical impact for both Segformer and EMSANet. However, as EMSANet has both better calibrated probability and uncertainty, we find it to be better for the object search task.

4.2 Perception's Impact on Sequential Navigation

Table 2. ObjectNav results of the Shortest-Path Policy with EMSANet.

| | Calibrated | Uncertainty SR Found | | Detection | | | |
|---------------------------|---------------|----------------------|------|-----------|------|------|------|
| Aggregation | Probabilities | Found Decision | | FPR | FNR | #FP | #FN |
| Ground Truth | × | × | 99.2 | 0.0 | 0.8 | 0.11 | 0.00 |
| Latest | × | × | 30.1 | 67.0 | 3.0 | 5.81 | 0.08 |
| Hits/Views | × | × | 69.2 | 12.4 | 18.4 | 0.52 | 0.19 |
| Skill Fusion [35] | × | × | 67.8 | 20.3 | 11.9 | 0.73 | 0.26 |
| Stubborn [17] | × | × | 32.7 | 63.8 | 3.5 | 5.49 | 0.08 |
| Latest Filtered | √ | × | 71.7 | 7.9 | 20.4 | 0.28 | 0.46 |
| Log odds | √ | \checkmark | 70.3 | 19.1 | 10.6 | 4.52 | 0.08 |
| Averaging | √ | × | 44.7 | 52.0 | 3.3 | 4.34 | 0.08 |
| Averaging | × | √ | 47.9 | 48.2 | 3.9 | 4.48 | 0.08 |
| Averaging (Ours) | √ | √ | 73.8 | 8.5 | 17.7 | 4.57 | 0.08 |
| Weighted Averaging (Ours) | √ | \checkmark | 74.9 | 8.7 | 16.4 | 4.58 | 0.08 |

Best and second best in bold and underlined. SR: success rate, FPR: false positive rate, NR: false positive rate, #FP: false positives, #FN: false negatives.

To isolate the impact of the perception in sequential navigation tasks, we first compare the perception-based aggregation strategies over identical sequences based on the shortest path policy. The results for the best performing semantic model, EMSANet, are reported in Tab. 2. We find that the shortest path policy achieves nearly perfect performance with the oracle ground truth perception, signifying any drop in the success rate can be attributed to the perception gap. Note that we find a very small number of six unsolvable episodes due to missing or incorrect targets. We find that temporal aggregation is essential, as simply using the latest predictions results in a high number of false found decisions, leading to a significant drop to 30%, while aggregation methods across the board reduce this gap. Secondly, we find both calibrated probabilities and uncertainty-based found decisions essential, outperforming the heuristic-based methods of using counts or erosion to eliminate false positives. Finally, we find our proposed averaging to provide the most reliable update. Qualitative examples are shown in Fig. 4. We can see that the agent is able to filter out the false positives (shown in circles), whenever the underlying map uncertainty is high for those objects.

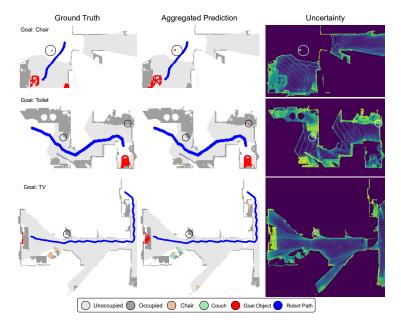


Fig. 4. Semantic maps showing, from left to right, the ground truth semantics, the aggregated predictions of our Weighted Averaging approach, and the resulting uncertainty map. Circles indicate positions where a target object was falsely detected but due to the high uncertainty, no false found decision was raised. The uncertainty varies from blue-yellow corresponding to 0.0-1.0 normalized entropy.

4.3 Perception's Impact on Sequential Decision Making

We then investigate the interdependence between the aggregation methods and policy, and evaluate the methods with the learned RL agent. The agent is trained with ground truth perception and deployed with the learned perception as described in Sec. 3.6. Results are presented in Tab. 3. The agent achieves a success rate of 75% with ground truth perception, indicating a policy gap of 25%. Across aggregation methods, we find that our results generalize to different policies, confirming the previous conclusions on the shortest path policy. Calibrated perception probabilities and uncertainties are essential which with our proposed aggregation method achieves a significant reduction in the perception gap, both in terms of success rate and path efficiency as measured by the SPL.

4.4 Generalizability Across Perception Models

Fig. 1 compares the results across the different perception models. We find the importance of calibrated aggregation confirmed across the board. While Mask-RCNN's object segmentation performs well in single image prediction, using newer models can close the often reported perception gap by over 10ppt with accurate aggregation. The full results across all aggregation methods, reported in Appx. A, confirm the relative results across aggregation methods, as the

Table 3. ObjectNav results of the RL policy trained with ground-truth semantics and deployed with EMSANet.

| | Calibrated | Uncertainty | \mathbf{SR} | Found | | Detection | | SPL |
|---------------------------|---------------|----------------|---------------|-------|------|-----------|------|------|
| Aggregation | Probabilities | Found Decision | | FPR | FNR | #FP | #FN | |
| Ground Truth | × | × | 75.0 | 3.8 | 21.2 | 0.18 | 0.00 | 27.9 |
| Latest | × | × | 15.8 | 81.8 | 2.4 | 3.51 | 0.13 | 7.1 |
| Hits/Views | × | × | 48.3 | 20.9 | 30.8 | 0.85 | 0.51 | 14.0 |
| Skill Fusion [35] | × | × | 48.4 | 33.1 | 18.5 | 0.89 | 0.46 | 16.7 |
| Stubborn [17] | × | × | 18.1 | 77.9 | 4.0 | 3.52 | 0.14 | 7.9 |
| Latest Filtered | √ | × | 51.2 | 20.0 | 28.8 | 0.44 | 0.70 | 13.9 |
| Log odds | √ | \checkmark | 52.9 | 27.1 | 20.0 | 5.26 | 0.24 | 19.2 |
| Averaging | √ | × | 28.5 | 61.2 | 10.3 | 3.36 | 0.17 | 10.7 |
| Averaging | × | \checkmark | 32.2 | 54.5 | 13.3 | 3.73 | 0.16 | 12.0 |
| Averaging (Ours) | √ | \checkmark | <u>55.6</u> | 15.3 | 29.1 | 5.73 | 0.19 | 19.2 |
| Weighted Averaging (Ours) | √ | \checkmark | 56.4 | 15.6 | 28.0 | 6.16 | 0.18 | 19.2 |

Best and second best in bold and underlined. SR: success rate, FPR: false positive rate, NR: false positive rate, #FP: false positives, #FN: false negatives.

relative ranking of methods remains very stable. We find the proposed weighted averaging to consistently perform best or on par across models and policies. One exception is the Skillfusion aggregation in combination with Segformer which it was developed with, evaluated with the RL policy. While achieving the same SPL, it reaches a 1.4ppt higher success rate in this setting. However, it is expected as we can see from the results in Fig. 3, that even after temperature scaling, the Segformer model's probabilities do not calibrate well. Nonetheless, utilizing calibrates uncertainties helps our method achieve a performance close to Skillfusion. This result further bolsters our claim that better calibrated networks help in better navigation tasks.

5 Conclusion

In this work, we identified gaps in the object search literature and showed that parts of the perception gap can be explained by the ineffective use of pretrained semantic perception models. We proposed the incorporation of calibrated probabilities and uncertainties across map aggregation and decision making, and demonstrated their effectiveness across different perception models and policies. The resulting methods are easy to incorporate and enable a wide family of models to further reduce the gap without additional training. As a result, they enable researchers to further close the gap. Furthermore, we made our code publicly available to aid future research.

In future work, we plan to further investigate direct policy conditioning on imperfect perception and uncertainties. To address the overconfidence induced by training RL agents with ground truth perception, we hypothesize that providing calibrated uncertainty maps as inputs can enable the agent to make better decisions such as investigating uncertain areas more closely and to learn to make context-dependent found decisions.

References

- 1. CVPR Embodied AI Workshop. https://embodied-ai.org/cvpr2023/ (2023)
- Bergstra, J., Bardenet, R., Bengio, Y., Kégl, B.: Algorithms for hyper-parameter optimization. In: Proc. of the Conf. on Neural Information Processing Systems, vol. 24 (2011)
- 3. Chang, M., Gervet, T., Khanna, M., Yenamandra, S., Shah, D., Min, S.Y., Shah, K., Paxton, C., Gupta, S., Batra, D., et al.: Goat: Go to any thing. arXiv preprint arXiv:2311.06430 (2023)
- 4. Chaplot, D.S., Gandhi, D., Gupta, S., Gupta, A., Salakhutdinov, R.: Learning to explore using active neural slam. In: Int. Conf. on Learning Representations (2020)
- Chaplot, D.S., Gandhi, D.P., Gupta, A., Salakhutdinov, R.R.: Object goal navigation using goal-oriented semantic exploration. Proc. of the Conf. on Neural Information Processing Systems 33, 4247–4258 (2020)
- 6. Chen, J., Li, G., Kumar, S., Ghanem, B., Yu, F.: How to not train your dragon: Training-free embodied object goal navigation with semantic frontiers. Robotics: Science and Systems (2023)
- Deitke, M., Han, W., Herrasti, A., Kembhavi, A., Kolve, E., et al.: Robothor: An open simulation-to-real embodied ai platform. In: Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (2020)
- Gal, Y., Ghahramani, Z.: Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In: Int. Conf. on Mach. Learning, pp. 1050–1059. PMLR (2016)
- Gosala, N., Petek, K., Drews-Jr, P.L., Burgard, W., Valada, A.: Skyeye: Self-supervised bird's-eye-view semantic mapping using monocular frontal view images. In: Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition, pp. 14,901–14,910 (2023)
- Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q.: On calibration of modern neural networks. In: Int. Conf. on Mach. Learning, pp. 1321–1330. PMLR (2017)
- 11. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Int. Conf. on Computer Vision, pp. 2961–2969 (2017)
- Honerkamp, D., Buchner, M., Despinoy, F., Welschehold, T., Valada, A.: Languagegrounded dynamic scene graphs for interactive object search with mobile manipulation. IEEE Robotics and Automation Letters (2024)
- 13. Hurtado, J.V., Valada, A.: Semantic scene segmentation for robotics. In: Deep learning for robot perception and cognition, pp. 279–311. Elsevier (2022)
- Jiang, J., Zheng, L., Luo, F., Zhang, Z.: Rednet: Residual encoder-decoder network for indoor rgb-d semantic segmentation. arXiv preprint arXiv:1806.01054 (2018)
- Kolve, E., Mottaghi, R., Han, W., VanderBilt, E., Weihs, L., Herrasti, A., Deitke, M., Ehsani, K., Gordon, D., Zhu, Y., et al.: Ai2-thor: An interactive 3d environment for visual ai. arXiv preprint arXiv:1712.05474 (2017)
- Lakshminarayanan, B., Pritzel, A., Blundell, C.: Simple and scalable predictive uncertainty estimation using deep ensembles. Proc. of the Conf. on Neural Information Processing Systems 30 (2017)
- 17. Luo, H., Yue, A., Hong, Z.W., Agrawal, P.: Stubborn: A strong baseline for indoor object navigation. In: Int. Conf. on Intelligent Robots and Systems, pp. 3287–3293. IEEE (2022)
- Maksymets, O., Cartillier, V., Gokaslan, A., Wijmans, E., Galuba, W., Lee, S., Batra, D.: Thda: Treasure hunt data augmentation for semantic navigation. In: Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition, pp. 15,374–15,383 (2021)

- 19. McCormac, J., Handa, A., Davison, A., Leutenegger, S.: Semanticfusion: Dense 3d semantic mapping with convolutional neural networks. In: Int. Conf. on Robotics & Automation, pp. 4628–4635 (2017)
- Naeini, M.P., Cooper, G., Hauskrecht, M.: Obtaining well calibrated probabilities using bayesian binning. Proc. of the National Conference on Artificial Intelligence 29(1) (2015)
- Nallapareddy, M.R., Sirohi, K., Drews, P.L., Burgard, W., Cheng, C.H., Valada, A.: Evcenternet: Uncertainty estimation for object detection using evidential learning. In: Int. Conf. on Intelligent Robots and Systems, pp. 5699–5706. IEEE (2023)
- 22. Ramakrishnan, S.K., Chaplot, D.S., Al-Halah, Z., Malik, J., Grauman, K.: Poni: Potential functions for objectgoal navigation with interaction-free learning. In: Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition. IEEE (2022)
- Ramakrishnan, S.K., Gokaslan, A., Wijmans, E., Maksymets, O., et al.: Habitat-matterport 3d dataset (HM3d): 1000 large-scale 3d environments for embodied AI.
 In: Proc. of the Conf. on Neural Information Processing Systems (2021)
- 24. Ramrakhya, R., Batra, D., Wijmans, E., Das, A.: Pirlnav: Pretraining with imitation and rl finetuning for objectnav. In: Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition, pp. 17,896–17,906 (2023)
- Raychaudhuri, S., Campari, T., Jain, U., Savva, M., Chang, A.X.: Mopa: Modular object navigation with pointgoal agents. In: Proc. of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 5763–5773 (2024)
- Roberts, M., Ramapuram, J., Ranjan, A., Kumar, A., Bautista, M.A., Paczan, N., Webb, R., Susskind, J.M.: Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In: Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition, pp. 10,912–10,922 (2021)
- Schmalstieg, F., Honerkamp, D., Welschehold, T., Valada, A.: Learning long-horizon robot exploration strategies for multi-object search in continuous action spaces. Proceedings of the International Symposium on Robotics Research (ISRR) (2022)
- 28. Schmalstieg, F., Honerkamp, D., Welschehold, T., Valada, A.: Learning hierarchical interactive multi-object search for mobile manipulation. IEEE Robotics and Automation Letters (2023)
- 29. Seichter, D., Fischedick, S.B., Köhler, M., Groß, H.M.: Efficient multi-task rgb-d scene analysis for indoor environments. In: 2022 Int. joint conference on neural networks (IJCNN), pp. 1–10. IEEE (2022)
- Sensoy, M., Kaplan, L., Kandemir, M.: Evidential deep learning to quantify classification uncertainty. Proc. of the Conf. on Neural Information Processing Systems 31 (2018)
- Sirohi, K., Büscher, D., Burgard, W.: uplam: Robust panoptic localization and mapping leveraging perception uncertainties. arXiv preprint arXiv:2402.05840 (2024)
- 32. Sirohi, K., Marvi, S., Büscher, D., Burgard, W.: Uncertainty-aware lidar panoptic segmentation. In: Int. Conf. on Robotics & Automation, pp. 8277–8283. IEEE (2023)
- Sirohi, K., Marvi, S., Büscher, D., Burgard, W.: Uncertainty-aware panoptic segmentation. IEEE Robotics and Automation Letters 8(5), 2629–2636 (2023)
- 34. Song, S., Lichtenberg, S.P., Xiao, J.: Sun rgb-d: A rgb-d scene understanding benchmark suite. In: Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition, pp. 567–576 (2015)
- 35. Staroverov, A., Muravyev, K., Yakovlev, K., Panov, A.I.: Skill fusion in hybrid robotic framework for visual object goal navigation. Robotics **12**(4), 104 (2023)
- 36. Szot, A., Clegg, A., Undersander, E., Wijmans, E., et al.: Habitat 2.0: Training home assistants to rearrange their habitat. arXiv preprint arXiv:2106.14405 (2021)

- 37. Thrun, S., Burgard, W.: Probabilistic robotics. Communications of the ACM 45(3), 52–57 (2002)
- 38. Wani, S., Patel, S., Jain, U., Chang, A.X., Savva, M.: Multi-on: Benchmarking semantic map memory using multi-object navigation. In: Proc. of the Conf. on Neural Information Processing Systems (2020)
- 39. Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J.M., Luo, P.: Segformer: Simple and efficient design for semantic segmentation with transformers. Proc. of the Conf. on Neural Information Processing Systems 34, 12,077–12,090 (2021)
- Yadav, K., Krantz, J., Ramrakhya, R., Ramakrishnan, S.K., Yang, J., et al.: Habitat challenge 2023. https://aihabitat.org/challenge/2023/ (2023)
- 41. Yamauchi, B.: A frontier-based approach for autonomous exploration. In: Proc. of the IEEE Int. Symp. on Comput. Intell. in Rob. and Aut. (1997)
- 42. Ye, J., Batra, D., Das, A., Wijmans, E.: Auxiliary tasks and exploration enable objectgoal navigation. In: Int. Conf. on Computer Vision, pp. 16,117–16,126 (2021)
- 43. Younes, A., Honerkamp, D., Welschehold, T., Valada, A.: Catch me if you hear me: Audio-visual navigation in complex unmapped environments with moving sounds. IEEE Robotics and Automation Letters 8(2), 928–935 (2023)
- Zemskova, T., Staroverov, A., Muravyev, K., Yudin, D., Panov, A.: Interactive semantic map representation for skill-based visual object navigation. IEEE Access (2024)
- 45. Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., Torralba, A.: Scene parsing through ade20k dataset. In: Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (2017)
- Zhou, K., Zheng, K., Pryor, C., Shen, Y., Jin, H., Getoor, L., Wang, X.E.: Esc: Exploration with soft commonsense constraints for zero-shot object navigation. arXiv preprint arXiv:2301.13166 (2023)
- Zhu, M., Zhao, B., Kong, T.: Navigating to objects in unseen environments by distance prediction. In: Int. Conf. on Intelligent Robots and Systems, pp. 10,571– 10,578. IEEE (2022)
- 48. Zürn, J., Weber, S., Burgard, W.: Trackletmapper: Ground surface segmentation and mapping from traffic participant trajectories. In: Proc. of the Conference on Robot Learning (2022)

A Appendix: Extended Results

We report full results with the shortest path policy for Mask-RCNN in Tab. 4 and Segformer in Tab. 5 and results using the RL Policy for both models in Tab. 6 and Tab. 7. Note that not all aggregation methods are applicable to Mask-RCNN.

Table 4. ObjectNav results of the Shortest-Path Policy with Mask-RCNN.

| | Calibrated | Uncertainty | SR | Found | | Dete | ction |
|-------------------|----------------|---------------|------|-------|------|------|-------|
| Aggregation | Found Decision | Probabilities | | FPR | FNR | #FP | #FN |
| Ground Truth | × | × | 99.2 | 0.00 | 0.8 | 0.11 | 0.00 |
| Latest | × | × | 46.1 | 32.2 | 21.7 | 2.11 | 0.39 |
| Hits/Views | × | × | 49.4 | 29.0 | 21.6 | 0.47 | 0.30 |
| Skill Fusion [35] | × | × | 59.2 | 13.7 | 27.1 | 0.32 | 0.55 |

Best and second best in bold and underlined. SR: success rate, FPR: false positive rate, NR: false positive rate, #FP: false positives, #FN: false negatives.

 ${\bf Table~5.~Object Nav~results~of~the~Shortest-Path~Policy~with~Segformer.}$

| | Calibrated | Uncertainty | SR | Found | | d Detect | |
|---------------------------|---------------|----------------|------|-------|------|----------|------|
| Aggregation | Probabilities | Found Decision | | FPR | FNR | #FP | #FN |
| Ground Truth | × | × | 99.2 | 0.00 | 0.8 | 0.11 | 0.00 |
| Latest | × | × | 41.9 | 46.8 | 11.3 | 3.04 | 0.22 |
| Hits/Views | × | × | 54.6 | 11.2 | 34.2 | 0.88 | 0.16 |
| Skill Fusion [35] | × | × | 66.0 | 11.7 | 22.3 | 0.26 | 0.47 |
| Stubborn [17] | × | × | 42.9 | 45.6 | 11.5 | 2.96 | 0.22 |
| Latest Filtered | √ | × | 64.6 | 11.8 | 23.6 | 0.27 | 0.57 |
| Log odds | √ | \checkmark | 63.6 | 19.9 | 16.5 | 2.18 | 0.23 |
| Averaging | √ | × | 53.1 | 34.9 | 12.0 | 2.21 | 0.23 |
| Averaging | × | \checkmark | 52.9 | 35.5 | 11.6 | 2.28 | 0.22 |
| Averaging (Ours) | √ | \checkmark | 65.8 | 6.5 | 27.7 | 2.23 | 0.23 |
| Weighted Averaging (Ours) | √ | \checkmark | 70.4 | 13.8 | 15.8 | 2.23 | 0.22 |

Best and second best in bold and underlined. SR: success rate, FPR: false positive rate, NR: false positive rate, #FP: false positives, #FN: false negatives.

Table 6. ObjectNav results of the RL policy trained with ground-truth semantics and deployed with Mask-RCNN.

| | Calibrated | Uncertainty | SR Found | | Dete | SPL | | |
|-------------------|---------------|----------------|----------|------|------|------|------|------|
| Aggregation | Probabilities | Found Decision | | FPR | FNR | #FP | #FN | |
| Ground Truth | × | × | 75.0 | 3.8 | 21.2 | 0.18 | 0.00 | 27.9 |
| Latest | × | × | 34.7 | 38.2 | 13.6 | 2.78 | 0.49 | 11.1 |
| Hits/Views | × | × | 40.2 | 34.3 | 25.5 | 1.11 | 0.44 | 12.2 |
| Skill Fusion [35] | × | × | 46.2 | 28.8 | 25.0 | 0.47 | 0.85 | 14.5 |

Best and second best in bold and underlined. SR: success rate, FPR: false positive rate, NR: false positive rate, #FP: false positives, #FN: false negatives.

Table 7. ObjectNav results of the RL policy trained with ground-truth semantics and deployed with Segformer.

| | Calibrated | Uncertainty | \mathbf{SR} | Found | | Detection | | SPL |
|---------------------------|---------------|----------------|---------------|-------|------|-----------|-----------------|------|
| Aggregation | Probabilities | Found Decision | | FPR | FNR | #FP | $\#\mathrm{FN}$ | |
| Ground Truth | × | × | 75.0 | 3.8 | 21.2 | 0.18 | 0.00 | 27.9 |
| Latest | × | × | 33.0 | 52.4 | 14.6 | 2.63 | 0.29 | 12.3 |
| Hits/Views | × | × | 50.4 | 15.8 | 33.8 | 1.04 | 0.29 | 12.9 |
| Skill Fusion [35] | × | × | 54.3 | 21.3 | 24.4 | 0.45 | 0.68 | 18.1 |
| Stubborn [17] | × | × | 36.7 | 52.9 | 10.4 | 2.58 | 0.27 | 14.0 |
| Latest Filtered | √ | × | 49.1 | 20.4 | 30.5 | 0.46 | 0.75 | 15.0 |
| Log odds | √ | \checkmark | 51.5 | 21.8 | 26.7 | 2.23 | 0.36 | 18.1 |
| Averaging | √ | × | 42.8 | 34.7 | 22.5 | 2.16 | 0.31 | 14.7 |
| Averaging | × | \checkmark | 44.9 | 33.2 | 21.9 | 2.25 | 0.31 | 16.2 |
| Averaging (Ours) | √ | \checkmark | 51.9 | 12.3 | 35.8 | 2.73 | 0.38 | 17.4 |
| Weighted Averaging (Ours) | √ | \checkmark | 52.9 | 17.6 | 29.5 | 2.44 | 0.36 | 18.1 |

Best and second best in bold and underlined. SR: success rate, FPR: false positive rate, NR: false positive rate, #FP: false positives, #FN: false negatives.