# Learning-based Relational Object Matching Across Views

Cathrin Elich[1,2], Iro Armeni[3], Martin R. Oswald[4], Marc Pollefeys[3,5], Joerg Stueckler[1]

arXiv:2305.02398v1 [cs.CV] 3 May 2023

*Abstract*— **Intelligent robots require object-level scene understanding to reason about possible tasks and interactions with the environment. Moreover, many perception tasks such as scene reconstruction, image retrieval, or place recognition can benefit from reasoning on the level of objects. While keypoint-based matching can yield strong results for finding correspondences for images with small to medium view point changes, for large view point changes, matching semantically on the object-level becomes advantageous. In this paper, we propose a learning-based approach which combines local keypoints with novel object-level features for matching object detections between RGB images. We train our object-level matching features based on appearance and inter-frame and cross-frame spatial relations between objects in an associative graph neural network. We demonstrate our approach in a large variety of views on realistically rendered synthetic images. Our approach compares favorably to previous state-of-the-art object-level matching approaches and achieves improved performance over a pure keypoint-based approach for large view-point changes.**

## I. Introduction

Object-centric representations are essential for robots that act in their environment, for instance, to perceive obstacles for navigation, or to plan object manipulation actions. Several approaches have been proposed that reason about scenes on the object-level for static and dynamic scene reconstruction [25], [31], [9]. When tracking multiple objects in image sequences, or recognizing if places are revisited for loop closure detection, the ability to match objects between views of the environment becomes important. Especially for large baselines changes between camera viewpoints, classical keypoint matching approaches [28] often break down due to drastic appearance changes and partial occlusions.

In this paper, we propose a novel approach for matching object detections between images by leveraging both spatial and appearance features of objects and spatial relations between the objects. In our approach, we combine keypoints with object-level feature correspondences to use the best of both worlds to match objects across small to large view point changes. On the object-level, the matching features are determined by an associative graph neural network that reasons on spatial and appearance features of object detections in both input images. To incorporate spatial knowledge into our matching pipeline, we additionally train our features on semantic and spatial auxiliary tasks such as inferring object class, 3D position, and relative distance between objects. By this, objects can be matched based on semantic and spatial information across large view point changes. This is in contrast to pure 2D keypoint-based matching approaches [26] which rely on local texture similarities and do not exploit object semantics. We evaluate our approach on a realistically rendered indoor dataset [21] and demonstrate improved accuracy over state-of-the-art object matching approaches. For large view-point changes, our approach also compares favorably in object matching to a baseline that uses the state-of-the-art keypoint matching method in accuracy and recall. We provide an ablation study to analyze the contributions of the individual design choices in our method.

In summary, we make the following contributions: **(1)** We present a novel learning approach for learning relational object features suited for matching detected objects between image pairs. Our approach uses an attentional graph neural network based on appearance and spatial features extracted from the object bounding boxes. **(2)** We combine our object-level features with learning-based keypoint matching to achieve state-of-the-art object matching for small to large view point changes.

## II. Related Work

Recently, several object-centric scene reconstruction approaches have been proposed that detect objects in images and aggregate the detections into scene representations like 3D object-level maps or scene graphs. For instance, methods for incremental simultaneous localization and mapping have been proposed which detect, segment, and reconstruct static and moving objects [17], [23], [35], [31], [24]. Scene graph-based approaches also estimate spatial and semantic relations between objects [1], [33], [22], [5]. Common to approaches which create object-centric scene representations from images is that they require means to detect or segment objects in the images and to associate them between images or instances in the accumulated scene representation. Typically this association is tackled using geometric data association.

In the multi-object tracking literature, data association is often performed using object-wise features extracted locally from the object bounding boxes which can aid in scenarios with occlusions [34]. Other video tracking approaches associate objects implicitly in neural encoder-decoder or transformer architectures [14], [2]. For scene graph fusion from objects detected in RGB images, CSR [4] learns encodings
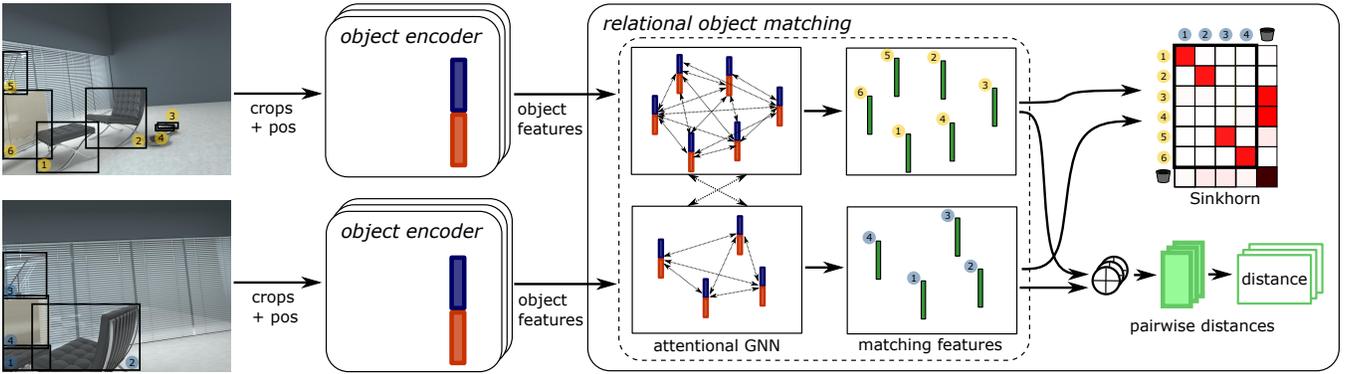
Fig. 1. **Learning relational object matching (ROM) features.** Our approach learns object-level features for matching objects between images. An object encoder network extracts object-wise features which encode appearance and positional information from the object bounding boxes. The features are input to an Attentional Graph Neural Network (AGNN) with self attention (within image) and cross attention (across images). The AGNN yields matching features for each object from which pairwise matching scores for objects of both images are computed by taking the dot product. An approximation to the optimal partial assignment between objects is determined using the Sinkhorn algorithm. The set of objects is augmented with a dustbin to allow for unmatched objects due to occlusions and limited field of view. Ground-truth object bounding boxes are used in the depicted example.

of object bounding boxes into feature vectors which are used to measure similarity of detected objects between views using cosine similarity. In Associative3D [19], also feature embeddings per objects are learned which are used to match objects between views. Both [27] and [30] leverage semantic and spatial information for finding 3D correspondences. Recently, [16] propose to use human models to estimate correspondences among wide-baseline view changes. Different to object associations, keypoint matching methods yield local image correspondences between pairs of images for moderate view point changes [26], [32], [10]. In our approach, inspired by SuperGlue [26], we further process individual object embeddings in a graph neural network which relates objects within images and across images to encode image context. Our approach can be an alternative or complementary approach for matching objects across views to geometric or keypoint-based approaches.

## III. METHOD

Our approach finds matches of corresponding objects in a 3D scene. Input to our method are pairs of images $I_1, I_2$ along with a set of object detections in each image. The objects can be extracted with any object detector such as YOLO [20], SSD [15], or Mask R-CNN [6]. Figs. 1 and 2 illustrate our approach. We extract deep encodings from each object detection. The encodings are further processed in an attentional graph neural network which establishes context within each image as well as across images. The resulting matching feature vectors are used to determine a similarity score for each possible pair of objects. Differentiable Sinkhorn matching [12], [3] then finds a unique best mapping for each object to either objects or an outlier class (dustbin). The network is trained end-to-end using ground truth matches of objects on realistically rendered images of indoor scenes. Finally, we propose an object matching pipeline which combines object-level features with learning-based keypoints (SuperGlue [26]) for matching in a wide range of view point changes.

### A. Object-Level Encoding

In each image, we encode the objects specified by 2D bounding boxes into two kinds of feature vectors which capture view-dependent and view-independent object information. We start from visual features $\mathbf{f}_{i,viz} \in \mathbb{R}^{D_{viz}}$ for each object $i$ which we determine by extracting ResNet34 [7] features pretrained on ImageNet within the object bounding boxes. The ResNet34 features are aggregated into a feature vector using a channel-wise max operation. Additionally, an MLP $g_{loc}$ calculates a feature vector $\mathbf{f}_{i,loc} \in \mathbb{R}^{D_{loc}}$ from the bounding box coordinates which we normalize by the image size. We use ReLU activation functions in each hidden layer. The visual and location features are concatenated into an input feature $\mathbf{f}_{i,in} := (\mathbf{f}_{i,viz}, \mathbf{f}_{i,loc})$. We use 3D regression and semantic object classification side-tasks to train features which are expressive in view-dependent and view-independent object properties. For the view-dependent features $\mathbf{f}_{i,dep} := g_{dep}(\mathbf{f}_{i,in})$, we apply an MLP on the input features and train a subsequent MLP $g_{pos}$ to predict the 3D position parametrized by a 2D image offset $\Delta\mathbf{p}_i \in \mathbb{R}^2$ from bounding box center to image projection of the 3D position and distance $d_i \in \mathbb{R}$, similar to Total3D [18]. The view-independent features $\mathbf{f}_{i,indep} := g_{indep}(\mathbf{f}_{i,in})$ are similarly extracted using an MLP, and another MLP $g_{class}$ extracts logits for semantic object classification. We obtain the final object feature $\mathbf{f}_i := (\mathbf{f}_{i,dep}, \mathbf{f}_{i,indep}) \in \mathbb{R}^{D_{obj}}$ by concatenating the view-dependent and view-independent features.

### B. Relational Object Matching

In a subsequent processing step, we use an Attentional Graph Neural Network (AGNN) as in [26] to obtain object encodings for matching. Object encodings are nodes in the graph, while graph edges model relations between objects within each image and across the images. The graph is fully connected and allows for interactions between all objects in both images. We denote object encodings by $\mathbf{x}_{k,i}$ where $k \in \{1, 2\}$ refers to the image of the object and $i$ specifies the object index. The object encodings
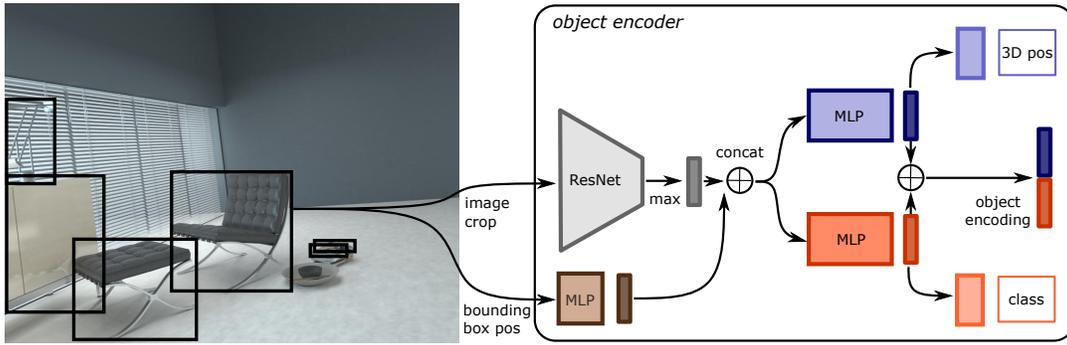
Fig. 2. **Object encoder.** From object bounding boxes, a ResNet34 pretrained on ImageNet extracts features which are max-aggregated channelwise into appearance feature vectors per object. The appearance features are concatenated with an MLP-encoding of the bounding box coordinates of the object. This object input feature is further processed by two MLPs into view-independent and view-dependent features. The first is trained with a 3D object position side task, while the latter is trained for semantic object classification. Both features are concatenated into an intermediate object feature.

$\mathbf{x}_{k,i,l} = \mathbf{x}_{k,i,l-1} + g_{AGNN}\big((\mathbf{x}_{k,i,l-1}, \mathbf{m}_{i.l})\big)$ are refined in each AGNN stage $l$ by passing message $\mathbf{m}_i$ aggregated from all other object nodes in the graph, i.e. within the same image and across images. The encodings are initialized with the respective object features $\mathbf{f}_i$. The messages are determined by an attention mechanism. An affine layer extracts values $\mathbf{v}_j$ and keys $\mathbf{k}_j$ for each other object in the graph from the current object encodings $\mathbf{x}_j$. Here, we drop the image index for simplicity of notation, since objects from both images are treated equally. For the object to be updated, we also determine the query $\mathbf{q}_i$ using an affine layer for its current encoding. The messages are computed by aggregating the values $\mathbf{m}_i = \sum_{j:(i,j)\in E} \alpha_{ij}\mathbf{v}_j$ with attention weights $\alpha_{ij} = \text{softmax}_j\left(\mathbf{q_i}^\top \mathbf{k}_j\right)$ and graph edges $E$.

The AGNN yields refined object features $\mathbf{x}_{k,i} \in \mathbb{R}^{D_{AGNN}}$ for object $i$ in image $k$ which are used to compute a similarity score using a scalar product between the feature vectors of objects in different images. In an additional 3D side task, we train the features to encode view-independent relative distance information between objects in each image. To this end, the relative distance $d_{k,ij} := g_{dist}\left((\mathbf{x}_{k,i}, \mathbf{x}_{k,j})\right)$ is regressed from pairs of refined features in the same image $k$ using an MLP. For further details on the AGNN, please refer to [26]. Finally, we use the Sinkhorn algorithm [12], [3] to compute a differentiable approximation to the optimal matching of objects based on their matching score. Let $M, N$ be the number of objects in image 1 and 2, respectively. Input to the Sinkhorn algorithm is the pairwise matching score of the objects determined by the scalar product of their feature vectors $\mathbf{S}_{ij}^{obj} = \mathbf{x}_{1,i}^\top \mathbf{x}_{2,j}$. This leads to a $M \times N$ matrix $\mathbf{S}^{obj}$. As noted in [26], using unnormalized feature vectors allows the network to learn prediction confidence implicitly.

Some objects might not be visible in both images due to occlusions or limited field-of-view. To allow for objects being not visible and matchable between frames for the matching, we additionally include an outlier association by a dustbin class with a single learnable parameter for each object. The augmented score matrix is denoted by $\overline{\mathbf{S}}^{obj} \in \mathbb{R}^{(M+1)\times(N+1)}$. Finally, let $\mathbf{1}_K \in \mathbb{R}^K$ be a vector of ones in each dimension. The optimal matching problem is to find

unique matchings for each objects to either other objects or the dustbins for each image. The dustbins are allowed to be associated multiple times. This linear assignment problem is

$$\arg\max_{\overline{\mathbf{P}}\in[0,1]^{M+1\times N+1}} \sum_{i,j} \overline{\mathbf{S}}_{ij}^{obj} \overline{\mathbf{P}}_{ij}$$

$$\text{s.t. } \overline{\mathbf{P}}\mathbf{1}_{N+1} = \left(\mathbf{1}_M^\top N\right)^\top \text{ and } \overline{\mathbf{P}}^\top \mathbf{1}_{M+1} = \left(\mathbf{1}_N^\top M\right)^\top \quad (1)$$

with assignment matrix $\overline{\mathbf{P}}$. The Sinkhorn algorithm iteratively normalizes rows and columns of $\exp\left(\overline{\mathbf{S}}^{obj}\right)$ to arrive at $\overline{\mathbf{P}}$. See [26] for additional details.

### C. Loss Function

We train our model supervised using ground truth object bounding boxes and matches between pairs of images. Our overall loss function

$$\mathcal{L} = \lambda_{aff}\mathcal{L}_{aff} + \lambda_{cls}\mathcal{L}_{cls} + \lambda_{pos}\mathcal{L}_{pos} + \lambda_{rel}\mathcal{L}_{rel} \quad (2)$$

is composed of four losses with corresponding weighting terms $\lambda_{aff}, \lambda_{cls}, \lambda_{pos}$ and $\lambda_{rel}$ accounting for affinity matching, classification, position, and relative distance. We use the affinity loss $\mathcal{L}_{aff}$ as in [26] which compares the predicted assignment matrix with the ground-truth assignment. The loss $\mathcal{L}_{cls}$ is a per-class weighted category cross entropy loss for the classification predicted from the view-independent features. The 3D position regression is trained with the loss

$$\mathcal{L}_{pos} = \frac{1}{M}\left(\sum_i (\mathbf{p}_{1,i} - \mathbf{p}_{1,i}^{gt})^2 + (d_{1,i} - d_{1,i}^{gt})^2\right)$$
$$+ \frac{1}{N}\left(\sum_j (\mathbf{p}_{2,j} - \mathbf{p}_{2,j}^{gt})^2 + (d_{2,i} - d_{2,i}^{gt})^2\right) \quad (3)$$

which measures the Euclidean distance between the predicted and the ground-truth position. We do not consider objects for this loss if the offset from the center of the bounding box is large (more than the bounding box length in each dimension), since these objects are typically found at the image border and lack image content for reliable prediction. The error in the prediction of the relative distance between

| MLP | hidden units, output dimensionality |
| --- | --- |
| $g_{loc}$ | [32, 64, 128] |
| $g_{dep/indep}$ | [512, 256, 128, 128] |
| $g_{pos}$ | [256, 3] |
| $g_{class}$ | [256, 40] |
| $g_{dist}$ | [256, 1] |
| $g_{AGNN}$ | [(self)256, (cross)256, (self)256, (cross)256] |

TABLE II

**MATCHING RESULTS FOR GROUND-TRUTH DETECTIONS AS INPUT. TOP: OBJECT-WISE, BOTTOM: FRAME-WISE. OUR APPROACH PERFORMS BEST AMONG TRAINED OBJECT MATCHING APPROACHES, AND OUTPERFORMS KEYPOINT-BASED MATCHING (SUPERGLUE) IN F1-SCORE AND RECALL FOR LARGE VIEWPOINT CHANGES.**

| object-wise | Easy | | | Hard | | | Very Hard | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | F1↑ | Prec↑ | Rec↑ | F1↑ | Prec↑ | Rec↑ | F1↑ | Prec↑ | Rec↑ |
| SuperGlue [26] | **0.657** | **0.811** | 0.552 | 0.353 | **0.682** | 0.250 | 0.169 | **0.574** | 0.099 |
| CSR [4] | 0.357 | 0.404 | 0.320 | 0.209 | 0.235 | 0.188 | 0.189 | 0.191 | 0.187 |
| Associative3D [19] | 0.259 | 0.265 | 0.253 | 0.191 | 0.180 | 0.204 | 0.194 | 0.164 | 0.239 |
| ROM features | 0.538 | 0.594 | 0.492 | 0.358 | 0.431 | 0.306 | 0.308 | 0.381 | 0.259 |
| Ours | 0.642 | 0.647 | **0.638** | **0.415** | 0.450 | **0.385** | 0.332 | 0.385 | **0.292** |
| **frame-wise** | F1↑ | Prec↑ | Rec↑ | F1↑ | Prec↑ | Rec↑ | F1↑ | Prec↑ | Rec ↑ |
| SuperGlue [26] | **0.740** | **0.811** | 0.589 | 0.551 | **0.671** | 0.278 | 0.489 | **0.557** | 0.099 |
| CSR [4] | 0.441 | 0.409 | 0.385 | 0.338 | 0.257 | 0.257 | 0.336 | 0.214 | 0.239 |
| Associative3D [19] | 0.328 | 0.312 | 0.312 | 0.274 | 0.228 | 0.264 | 0.284 | 0.204 | 0.282 |
| ROM features | 0.625 | 0.628 | 0.553 | 0.493 | 0.473 | 0.381 | 0.469 | 0.402 | 0.326 |
| Ours | 0.723 | 0.688 | **0.694** | 0.552 | 0.500 | **0.464** | 0.498 | 0.419 | **0.367** |

objects is measured using the loss

$$\mathcal{L}_{rel} = \sum_{k=1}^{2} \sum_{i} \sum_{j, j\neq i} (d_{k,ij} - d_{k,ij}^{gt})^2. \qquad (4)$$

### D. Combined Object-Level and Keypoint-based Matching

Classically, keypoints with local descriptors have been used to find correspondences between images. While modern learning-based approaches such as SuperGlue [26] can achieve highly accurate results, finding correspondences across large view points still poses a challenge for these methods, since they do not take object-level information into account. We thus propose to combine keypoint-based matching with our object-level feature matching. To this end, we find keypoint matches between two images using SuperGlue as state-of-the-art learning-based keypoint detector. Each keypoint match is assigned to those object matches for which the keypoints are located within the bounding box of the objects. The score $\mathbf{S}_{ij}^{kp}$ of each object match for object $i, j$ is determined by the logarithm of the count of keypoint matches between the objects. An additional dustbin is added for each object with a score of 1. The combined object matching score $\overline{\mathbf{S}} = \overline{\mathbf{S}}^{obj} + \alpha \overline{\mathbf{S}}^{kp}$ is a linear combination of the object-level and keypoint-based scores, from which we find the optimal assignment using the Sinkhorn algorithm. We use $\alpha = 100$ in our experiments which we chose based on the F1-score on the validation set.

## IV. EXPERIMENTS

We evaluate our method on realistically rendered indoor scenes with ground-truth object matchings. Our experiments address three key questions: **(1)** What is the matching performance of our approach for objects for various difficulty levels wrt. view point change? **(2)** How does our approach relate to state-of-the-art baselines based on object and keypoint matching? **(3)** What are the benefits of our design choices of auxiliary losses and differentiable matching?

**Datasets.** We evaluate our approach on the photorealistic synthetic Hypersim dataset of indoor scenes [21]. We filter scenes with unrealistic appearance as well as scenes or images containing less than three objects. We discard structural objects (e.g. wall, floor) as well as instances of the "otherprop" class and objects with minimum bounding box side length of 25 pixels. Structural objects are typically not well localized in position, whereas the objects in the "otherprop" class are typically small in the image with only little context within the bounding box for object feature extraction. We consider pairs of images from the same scene where at least two objects are depicted in both of them. For our experiments, we use the official train/val/test split which results in a final number of 302/43/40 rooms and a total number of approximately 45.8k/6.3k/6.3k images. Per training epoch, we iterate over all frames and sample a valid counterpart. For testing, we generated three subsets of image pairs with different level of difficulty by sampling one pair per frame in the test split in each category (if available) We consider (a) the average difference of distances $\overline{d}$ of objects from the respective cameras as well as (b) the average over angles $\overline{\alpha}$ between the object-to-camera rays. We denote the subsets as easy ($\overline{d} \leq 4\,m$, $\overline{\alpha} \leq 45°$), hard (remaining objects with $\overline{d} \leq 8\,m, \overline{\alpha} \leq 90°$), and very hard ($\overline{d} > 8\,m$ or $\overline{\alpha} > 90°$). The subsets consist of 45.7k / 6.2k / 6.3k (easy), 42.7k / 5.8k / 6.1k (hard), and 24.4k / 3.6k / 3.0k (very hard) image pairs for training / validation / testing, respectively. For the ablation and comparison to baselines, we provide results with GT object bounding boxes as input as a base result for an optimal detector. We also show results on predicted object detections from EfficientDet-D7 [36], [8] trained on the MS COCO dataset [13].

**Network Parameters and Training Details.** We apply a pretrained ResNet34 [7] to obtain object-wise visual features from 128×128 patches obtained by rescaling the detected object bounding boxes. Tab. I lists the network parameters of object encoder and matching network. The AGNN alternates self- and cross-attention 2 times. During training, we add Gaussian noise with zero mean and variance 0.01 to the precomputed visual feature vectors for data augmentation. We only consider a maximum number of 40 objects per frame and sample objects if more are present. Input bounding box coordinates are normalized wrt. image size to values between 0 and 1. The loss weights are set to $\{\lambda_{aff}, \lambda_{cls}, \lambda_{pos}, \lambda_{rel}\} = \{1, 1, 0.1, 0.1\}$. We found 10 Sinkhorn iterations to be sufficient to converge. We train our model using ADAM [11] with learning rate $10^{-4}$ and batch size 32 for 350 epochs.

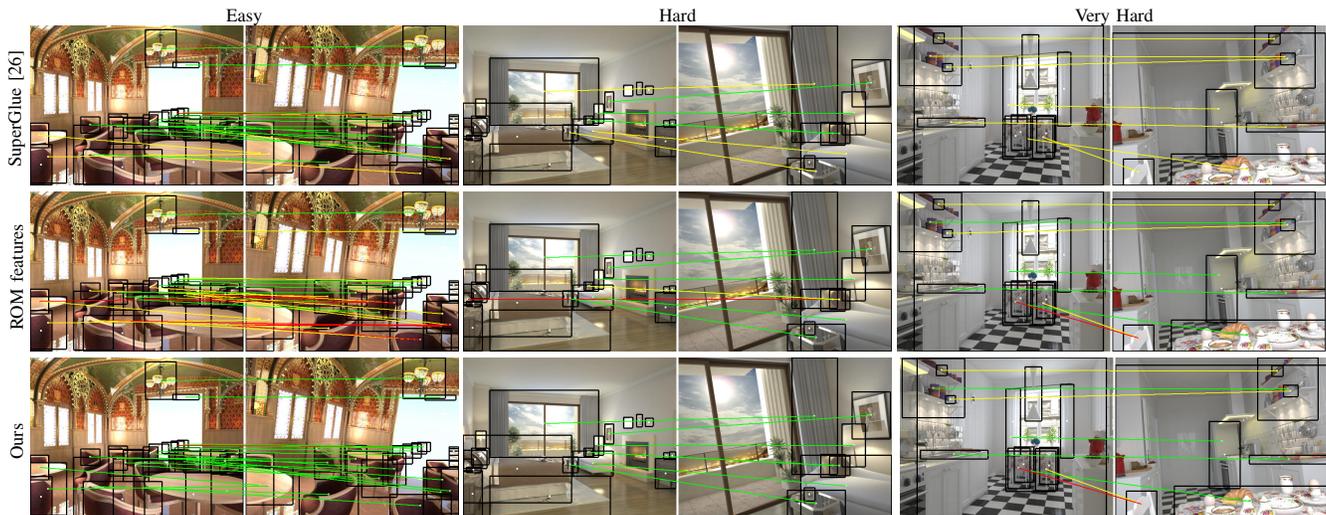| | Classification | 3D Position | | | | 3D Distance | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Acc ↑ | Mean ↓ | Median ↓ | Err($\leq 0.5$ m) ↑ | Err($\leq 1$ m) ↑ | Mean ↓ | Median ↓ | Err($\leq 0.5$ m) ↑ | Err($\leq 1$ m) ↑ |
| Associative3D [19] | **0.598** | 2.438 | 1.829 | 0.066 | 0.240 | 1.500 | **0.956** | 0.294 | **0.516** |
| Ours (gt bb) | 0.380 | 2.492 | **1.775** | **0.135** | **0.292** | **1.225** | 0.995 | **0.350** | 0.504 |



Fig. 3. **Qualitative comparison.** Evaluation on GT bounding boxes; **green:** correctly detected match, **red:** wrongly detected match, **yellow:** missing match. While SuperGlue yields very accurate matching results for the examples of easy and hard view point changes, it cannot detect any matches if the change of view point is too large (right column, very hard). On the other hand, our ROM-features can also be used to match objects in images with larger changes. Combining both object- and keypoint-based matching benefits from both capabilities.
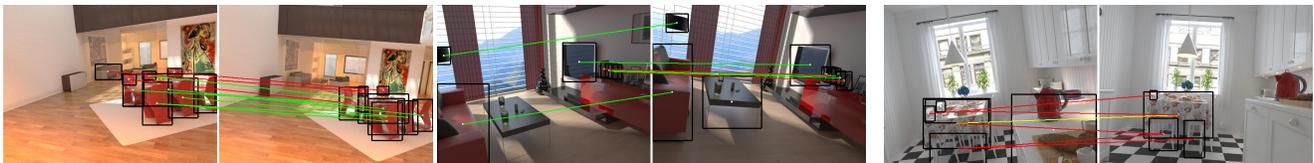


Fig. 4. **Qualitative results from our method.** Evaluation on detected bounding boxes; **green:** correctly detected match, **red:** wrongly detected match, **yellow:** missing match. Right: failure case with inaccurately detected object bounding boxes.

**Evaluation Metrics.** The main task of our approach is to obtain valid object matches in pairs of images. Given the assignment matrix $\overline{\mathbf{P}}$, we determine a fixed assignment per object for both images either to some object of the other image or to the dustbin by maximum value assignment. We compute the recall of correct object-object matches over all ground truth matches. Precision is defined as the ratio of correct object-object matches among all found object-object matches. We either compute the respective measure over all objects in all images combined (*object-wise*) or individually over all objects per frame-pair with subsequent averaging over all images (*frame-wise*). Please note that scenes with a very large number of objects have a high impact on the first-named variant. When evaluating our approach using objects from a detector network as input, the detector might detect different objects than actually available in the ground truth. Hence, we only evaluate the matching for detections which are assigned to a ground truth object in the same image. A

detection is assigned to a ground truth object if it overlaps and has maximal intersection-over-union with the ground truth bounding box. We further evaluate the performance of our models on the auxiliary tasks. For the classification task, we compute the accuracy of the semantic predictions over all object. Similar to [19], we compute mean and median errors, and the rate of estimates with errors $\leq 0.5/1.0\,m$ for both the distance and position prediction.

**Baselines.** We compare our method with two state-of-the-art baselines, CSR [4] and Associative3D [19], which learn object matching from bounding boxes. We trained the approaches on our dataset using the public reference implementation until convergence on the validation set. We further use SuperGlue [26] keypoints to determine an object-level matching by applying Sinkhorn iterations on $\mathbf{S}^{kp}$.

*A. Object Matching Results*

**Ground-truth detections.** In Tab. II we show matching results of our approach using ground truth object bounding

| | Easy | | | Hard | | | Very Hard | | |
|---|---|---|---|---|---|---|---|---|---|
| | F1↑ | Prec↑ | Rec↑ | F1↑ | Prec↑ | Rec↑ | F1↑ | Prec↑ | Rec↑ |
| only ResNet34 feats | 0.387 | 0.365 | 0.411 | 0.235 | 0.201 | 0.281 | 0.190 | 0.149 | 0.260 |
| w/o viz-input | 0.396 | **0.620** | 0.291 | 0.204 | 0.389 | 0.138 | 0.178 | 0.332 | 0.122 |
| w/o pos-input | 0.437 | 0.545 | 0.365 | 0.285 | 0.402 | 0.221 | 0.262 | 0.376 | 0.201 |
| w/o AGNN | 0.456 | 0.590 | 0.372 | 0.294 | 0.411 | 0.229 | 0.272 | 0.374 | 0.214 |
| only self-att. GNN | 0.517 | 0.613 | 0.446 | 0.344 | <u>0.452</u> | 0.277 | 0.292 | 0.399 | 0.231 |
| w/o $L_{cls}$ | 0.525 | <u>0.619</u> | 0.456 | 0.338 | **0.456** | 0.268 | 0.291 | **0.433** | 0.219 |
| w/o $L_{pos}$ | <u>0.533</u> | 0.589 | <u>0.487</u> | <u>0.357</u> | 0.428 | **0.306** | **0.315** | 0.379 | <u>0.269</u> |
| w/o $L_{dist}$ | 0.471 | 0.478 | 0.463 | 0.322 | 0.348 | 0.299 | 0.272 | 0.242 | **0.310** |
| w/o $L_{cls/pos/dist}$ | 0.486 | 0.574 | 0.421 | 0.330 | 0.433 | 0.267 | <u>0.310</u> | <u>0.403</u> | 0.252 |
| w/o $L_{aff}$ | 0.378 | 0.386 | 0.370 | 0.202 | 0.199 | 0.205 | 0.164 | 0.148 | 0.184 |
| ROM features (gt bb) | **0.538** | 0.594 | **0.492** | **0.358** | 0.431 | **0.306** | 0.308 | 0.381 | 0.259 |

| | all GT matchings | | | detections only | |
|---|---|---|---|---|---|
| | F1↑ | Prec↑ | Rec↑ | F1↑ | Rec↑ |
| Ours (gt bb) | 0.526 | 0.552 | 0.502 | - | - |
| SuperGlue [26] | 0.139 | **0.413** | 0.084 | **0.385** | 0.360 |
| ROM features | 0.095 | 0.168 | 0.066 | 0.212 | 0.285 |
| Ours | **0.143** | 0.230 | **0.103** | 0.303 | **0.445** |

boxes as input. We also compare our method with state-of-the-art approaches which are trained end-to-end for object description and matching (CSR, Associative3D) and a keypoint-based baseline (SuperGlue). It can be observed that both using only our ROM features as well as our combined method consistently outperforms CSR and Associative3D. SuperGlue has an advantage on smaller to medium (easy) view point changes, while our ROM features are on par (hard) or better (very hard) on larger view point changes in recall and F1 score. Especially, we found it challenging for SuperGlue to detect any keypoint matches in case of larger view point changes, e.g. when the same scene was depicted from opposite sides. By combining the two matching principles, we can benefit from the advantages of both methods. In most cases, inference and matching time for our ROM features took approx. $50 - 100$ ms per image pair. We show qualitative examples and comparisons with baselines in Fig. 3. In Tab. III, we further evaluate the performance of our network on the auxiliary tasks of classification, 3D position, and distance estimation and compare with Associative3D. Our approach estimates 3D position and distance of objects with a median error of $1.775\,\text{m}$ and $0.995\,\text{m}$, respectively, from RGB images.

Tab. IV shows results of various ablations of our model using ground truth object detections as input. Each architectural design choice contributes to the overall performance. The classification and 3D losses do only bring little improvements on the easy and hard cases, while the distance loss helps to improve the matching performance on all difficulty levels. Also both self- and cross-attention and using visual (viz-input) or position (pos-input) features are important for the performance of the method.

**Predicted object detections.** We also demonstrate our approach for predicted object detections in Tab. V. From the detections, we remove those whose predicted class label (COCO labels) does not appear in the NYU40 [29] class labels used in our GT. We further only consider detections with a minimum detection score of 0.3. Predicted bounding boxes are assigned to an associated GT object if the intersection over union (IoU) between GT and prediction is greater than 0.5. If multiple objects would get assigned to the same objects, we choose the one with highest IoU. By this, approx. 40% of the detections are assigned to GT objects and 30% of all GT objects are found. Contrary to using GT bounding boxes, less than three objects might get found in an image and forwarded to the model. We report F1 and recall scores for recovering the matching for the detected objects that have GT correspondence. We also show results for recovering the complete GT matching which, however, depends on the detector performance too. For matching available predicted detections, our approach achieves moderately lower scores than for matching GT detections.

### B. Limitations and Future Work

Our matching approach relies on a pretrained detector and the quality of the detected bounding boxes. In future work, end-to-end training of the detector with the matching pipeline could be investigated. In difficult occlusion settings, object bounding boxes can overlap and similar features are computed. Calculating object features based on instance segmentation could be interesting for future research. In this work, we focused on photorealistic synthetic, static indoor scenes with full 3D object-wise position information during training which allows our full model to consider distances between objects. While our model also achieves decent performance without considering spatial information of objects, we believe that ideas for matching objects in case of scene changes (e.g. outdoor scenes with moving cars), different camera geometries, or limited supervision are further interesting directions to explore.

### V. CONCLUSION

We proposed a novel approach for matching 2D object detections between images. Our approach is trained to match objects across a large variety of view points with auxiliary classification and 3D regression tasks. This way, our approach can outperform keypoint-based matching approaches for large view point changes. We evaluate our approach using a realistically rendered indoor dataset, and also demonstrate state-of-the-art performance among approaches which train object-wise matching end-to-end. Future applications of our approach could be explored for tasks such as object-centric scene reconstruction, image retrieval, and localization.

REFERENCES

[1] I. Armeni, Z.-Y. He, J. Gwak, A. R. Zamir, M. Fischer, J. Malik, and S. Savarese, "3D Scene Graph: A structure for unified semantics, 3D space, and camera," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 5664–5673.

[2] Y. Cong, W. Liao, H. Ackermann, B. Rosenhahn, and M. Y. Yang, "Spatial-temporal transformer for dynamic scene graph generation," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 16 352–16 362.

[3] M. Cuturi, "Sinkhorn distances: Lightspeed computation of optimal transport," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2013, pp. 2292–2300.

[4] S. Y. Gadre, K. Ehsani, S. Song, and R. Mottaghi, "Continuous scene representations for embodied AI," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[5] N. Gothoskar, M. Cusumano-Towner, B. Zinberg, M. Ghavamizadeh, F. Pollok, A. Garrett, J. Tenenbaum, D. Gutfreund, and V. Mansinghka, "3DP3: 3D scene perception via probabilistic programming," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2021, pp. 9600–9612.

[6] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask R-CNN," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2017, pp. 2980–2988.

[7] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.

[8] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, and K. Murphy, "Speed/accuracy trade-offs for modern convolutional object detectors," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 3296–3297.

[9] N. Hughes, Y. Chang, and L. Carlone, "Hydra: A real-time spatial perception system for 3D scene graph construction and optimization," in *Robotics: Science and Systems (RSS)*, 2022.

[10] W. Jiang, E. Trulls, J. Hosang, A. Tagliasacchi, and K. M. Yi, "COTR: Correspondence transformer for matching across images," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.

[11] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *International Conference on Learning Representations (ICLR)*, 2014.

[12] P. Knopp and R. Sinkhorn, "Concerning nonnegative matrices and doubly stochastic matrices," *Pacific Journal of Mathematics*, pp. 343 – 348, 1967.

[13] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *European Conference on Computer Vision (ECCV)*. Springer International Publishing, 2014, pp. 740–755.

[14] C. Liu, Y. Jin, K. Xu, G. Gong, and Y. Mu, "Beyond short-term snippet: Video relation detection with spatio-temporal global context," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 10 837–10 846.

[15] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot multibox detector," in *European Conference on Computer Vision (ECCV)*. Springer, 2016, pp. 21–37.

[16] W.-C. Ma, A. J. Yang, S. Wang, R. Urtasun, and A. Torralba, "Virtual correspondence: Humans as a cue for extreme-view geometry," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 15 924–15 934.

[17] J. McCormac, R. Clark, M. Bloesch, A. J. Davison, and S. Leutenegger, "Fusion++: Volumetric object-level SLAM," in *International Conference on 3D Vision (3DV)*, 2018, pp. 32–41.

[18] Y. Nie, X. Han, S. Guo, Y. Zheng, J. Chang, and J. J. Zhang, "Total3DUnderstanding: Joint layout, object pose and mesh reconstruction for indoor scenes from a single image," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[19] S. Qian, L. Jin, and D. F. Fouhey, "Associative3D: Volumetric reconstruction from sparse views," in *European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 140–157.

[20] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, real-time object detection," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 779–788.

[21] M. Roberts, J. Ramapuram, A. Ranjan, A. Kumar, M. A. Bautista, N. Paczan, R. Webb, and J. M. Susskind, "Hypersim: A photo-realistic synthetic dataset for holistic indoor scene understanding," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.

[22] A. Rosinol, A. Violette, M. Abate, N. Hughes, Y. Chang, J. Shi, A. Gupta, and L. Carlone, "Kimera: From slam to spatial perception with 3D dynamic scene graphs," *The International Journal of Robotics Research*, pp. 1510–1546, 2021.

[23] M. Runz, M. Buffier, and L. Agapito, "MaskFusion: Real-time recognition, tracking and reconstruction of multiple moving objects," in *IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, 2018, pp. 10–20.

[24] M. Rünz, K. Li, M. Tang, L. Ma, C. Kong, T. Schmidt, I. D. Reid, L. Agapito, J. Straub, S. Lovegrove, and R. A. Newcombe, "FroDO: From detections to 3D objects," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 14 708–14 717.

[25] R. F. Salas-Moreno, R. A. Newcombe, H. Strasdat, P. H. Kelly, and A. J. Davison, "SLAM++: Simultaneous localisation and mapping at the level of objects," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013, pp. 1352–1359.

[26] P.-E. Sarlin, D. DeTone, T. Malisiewicz, and A. Rabinovich, "SuperGlue: Learning feature matching with graph neural networks," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[27] J. L. Schönberger, M. Pollefeys, A. Geiger, and T. Sattler, "Semantic visual localization," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 6896–6906.

[28] J. L. Schönberger and J.-M. Frahm, "Structure-from-motion revisited," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[29] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from RGBD images," in *European Conference on Computer Vision (ECCV)*. Springer, 2012, pp. 746–760.

[30] P. Speciale, D. P. Paudel, M. R. Oswald, H. Riemenschneider, L. Van Gool, and M. Pollefeys, "Consensus maximization for semantic region correspondences," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7317–7326.

[31] M. Strecke and J. Stückler, "EM-Fusion: Dynamic object-level SLAM with probabilistic data association," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 5864–5873.

[32] J. Sun, Z. Shen, Y. Wang, H. Bao, and X. Zhou, "LoFTR: Detector-free local feature matching with transformers," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

[33] J. Wald, H. Dhamo, N. Navab, and F. Tombari, "Learning 3D semantic scene graphs from 3D indoor reconstructions," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[34] N. Wojke, A. Bewley, and D. Paulus, "Simple online and realtime tracking with a deep association metric," in *IEEE International Conference on Image Processing (ICIP)*, 2017, pp. 3645–3649.

[35] B. Xu, W. Li, D. Tzoumanikas, M. Bloesch, A. J. Davison, and S. Leutenegger, "MID-Fusion: Octree-based object-level multi-instance dynamic SLAM," in *International Conference on Robotics and Automation (ICRA)*, 2019, pp. 5231–5237.

[36] H. Yu, C. Chen, X. Du, Y. Li, A. Rashwan, L. Hou, P. Jin, F. Yang, F. Liu, J. Kim, and J. Li, "TensorFlow Model Garden," https://github.com/tensorflow/models, 2020.