

VII. SUPPLEMENTARY MATERIAL

A. PR AUC and Recall@1 Per Sequence Pairing

To provide more insights into how the various methods reported in Table III perform on each of the sequence pairings, Fig. 5 shows the precision-recall curves resulting from exhaustive pairwise matching of each sequence combination. Since Multi-Process Fusion does not provide an interface for exhaustive matching, we are not able to add it to the precision-recall plots. We further report retrieval performance for $N = 1$ on each combination of query and database sequences in Table VII. As outlined in Section V-B, we compare our composite descriptor to both visual and structural descriptors. We further evaluate against NetVLAD variants obtained using either off-the-shelf (ots) or fine-tuned (ft) weights as well as against Multi-Process Fusion, DenseVLAD and SeqSLAM for sequence lengths of 20 m and 40 m, respectively. For this comparison, all descriptors including those produced by NetVLAD and DenseVLAD are projected to $dim_f = 256$ dimensions. Our composite descriptor outperforms the single-modality descriptors as well as all of the baselines in exhaustive pairwise matching of each sequence combination. On top of that, it shows best retrieval performance on all except two sequence pairings. Our results further reveal that fine-tuning of NetVLAD on our training dataset somewhat deteriorates the pairwise matching accuracy for several sequence combinations but improves retrieval performance for the majority of sequence pairings. Similar observations have been reported in [23].

B. Network Configurations

TABLE V

THE NUMBER OF OUTPUT CHANNELS AND THE ARRANGEMENT OF POOLING OPERATIONS FOR EACH OF THE d_S LAYERS OF THE STRUCTURE-BASED FEATURE EXTRACTION NETWORKS USED TO OBTAIN THE RESULTS REPORTED IN TABLE II.

d_S	# Channels per CONV-layer	Pooling indices
6	$2 \times [32], 2 \times [64], 2 \times [128]$	2, 3, 4
8	$2 \times [32], 3 \times [64], 3 \times [128]$	2, 4, 6
9	$2 \times [32], 4 \times [64], 3 \times [128]$	2, 4, 6, 8
10	$2 \times [32], 4 \times [64], 4 \times [128]$	2, 4, 6, 8
12	$2 \times [32], 5 \times [64], 5 \times [128]$	2, 4, 7, 10

In Table V, we provide an overview of all structure-based feature extraction network configurations evaluated and compared in Section V-A. The pooling indices denote which convolutional layers are followed by pooling operations – with the networks with depths $d_S = 6$ and $d_S = 12$ not following the pattern described in Section III-A. Furthermore, Table VI shows the feature extraction configurations used in our best performing compound network as well as the slightly modified versions used to obtain the 256-dimensional appearance and structural descriptors reported in Table III.

C. Image Retrieval

We provide additional examples of top retrieved database instances when using either appearance, structure, or com-

TABLE VI

OVERVIEW OF THE APPEARANCE- (A) AND STRUCTURE-BASED (S) FEATURE EXTRACTION NETWORK CONFIGURATIONS REPORTED IN TABLE III. LISTED ARE THE CONFIGURATION OF OUR BEST PERFORMING COMPOUND NETWORK AS WELL AS THE SLIGHTLY MODIFIED VERSIONS USED TO GENERATE THE 256-DIMENSIONAL SINGLE-MODALITY DESCRIPTORS.

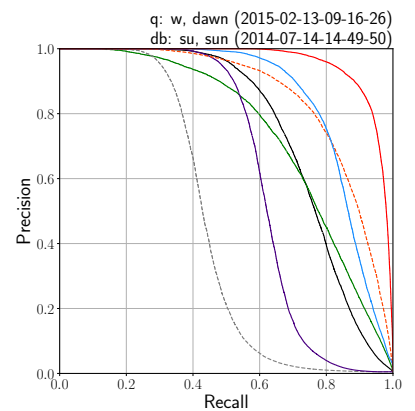
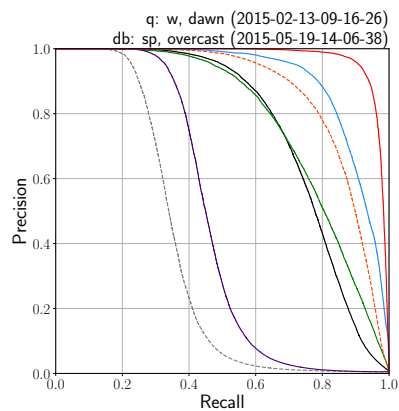
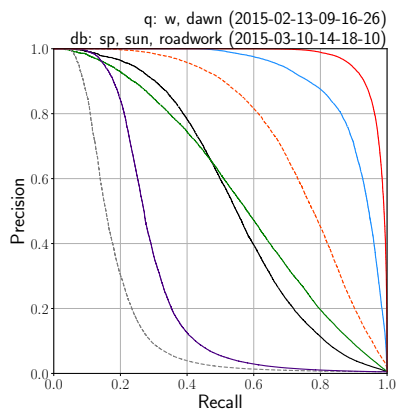
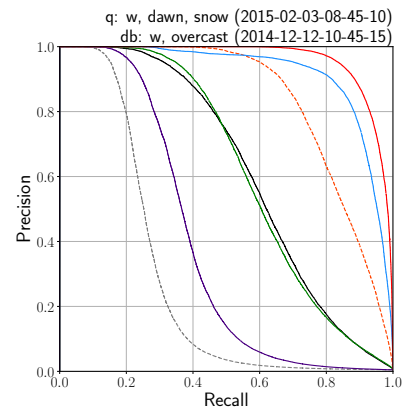
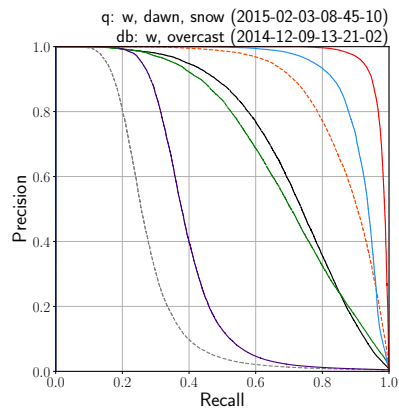
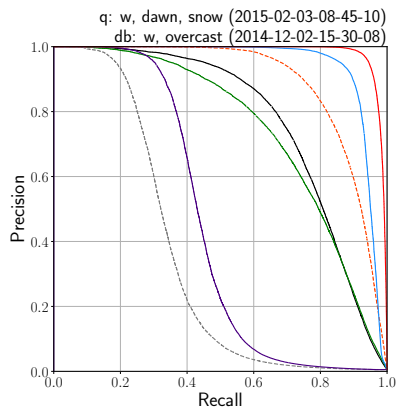
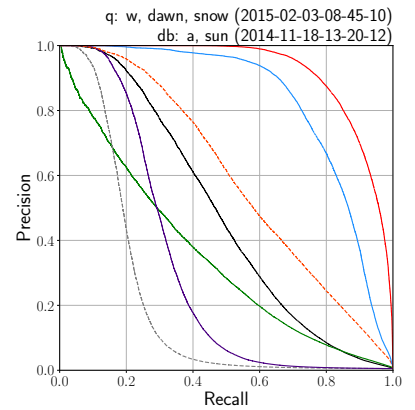
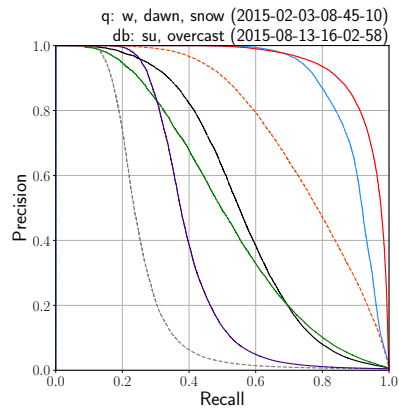
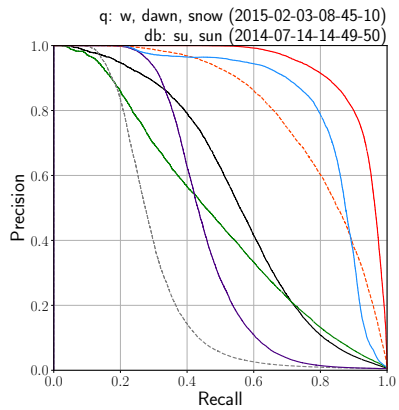
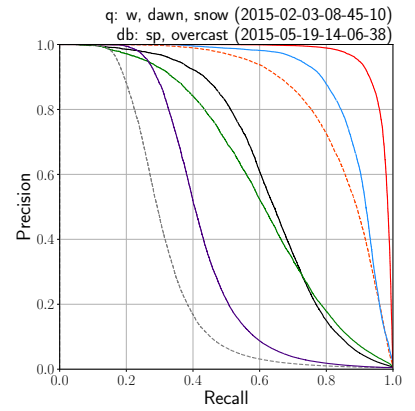
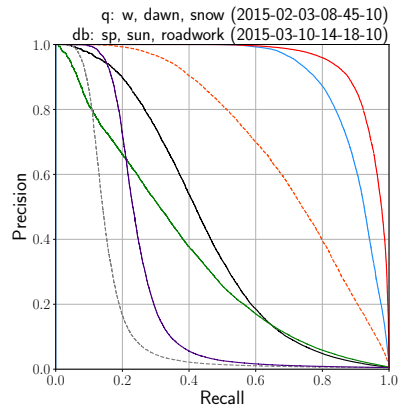
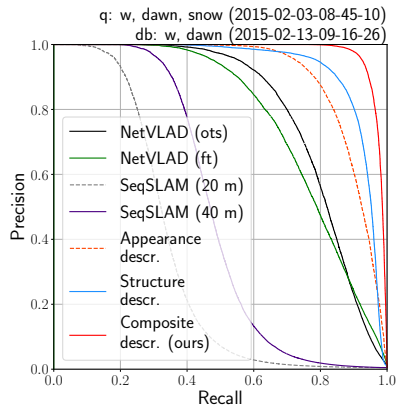
Type	d	# Channels per CONV-layer	Pooling indices
Compound network			
A	12	$6 \times [64], 6 \times [128]$	2, 4, 6, 8, 10
S	9	$2 \times [32], 4 \times [64], 3 \times [128]$	2, 4, 6, 8
Appearance- / structure-based networks			
A	12	$5 \times [64], 5 \times [128], 2 \times [256]$	2, 4, 6, 8, 10
S	9	$2 \times [32], 3 \times [64], 2 \times [128], 2 \times [256]$	2, 4, 6, 8

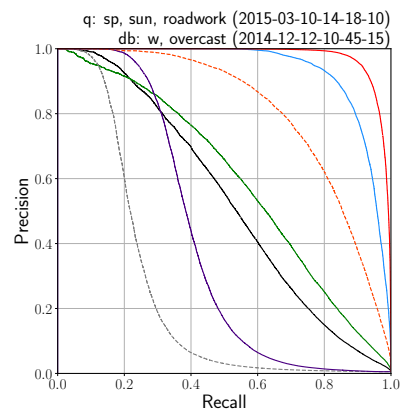
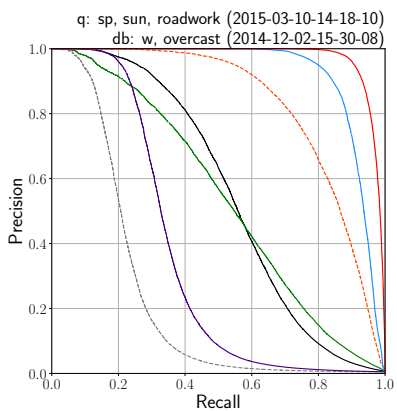
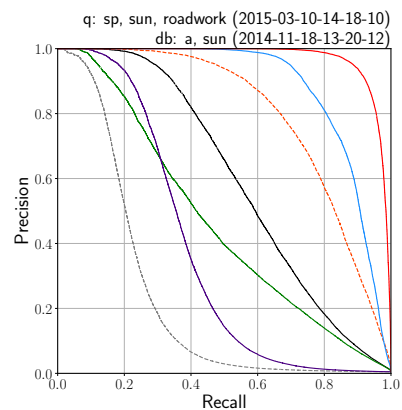
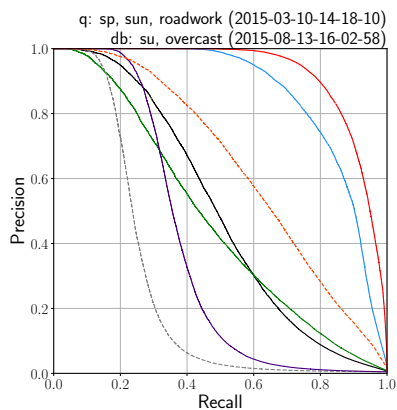
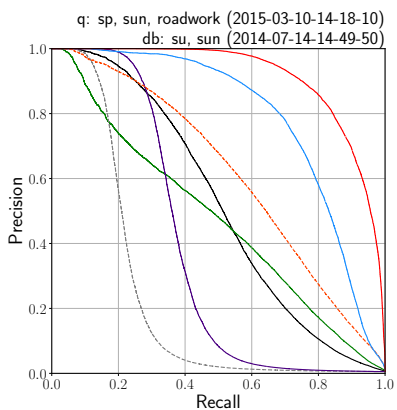
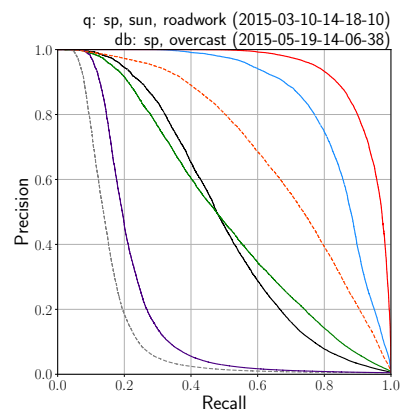
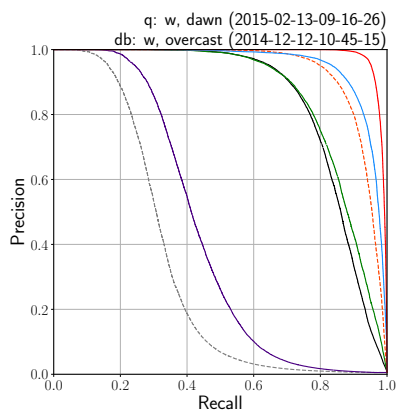
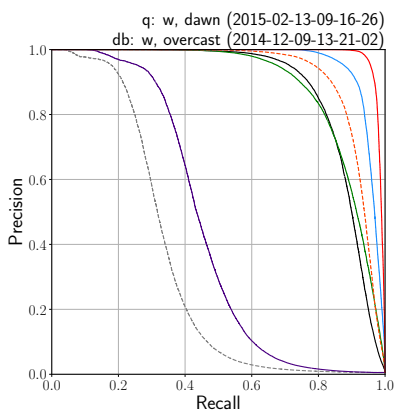
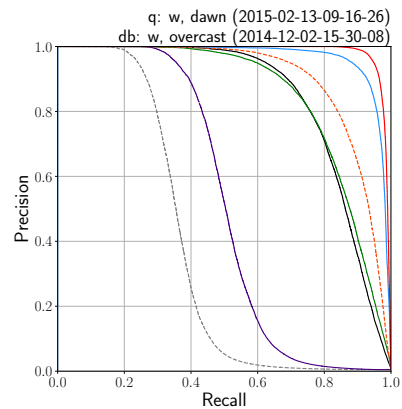
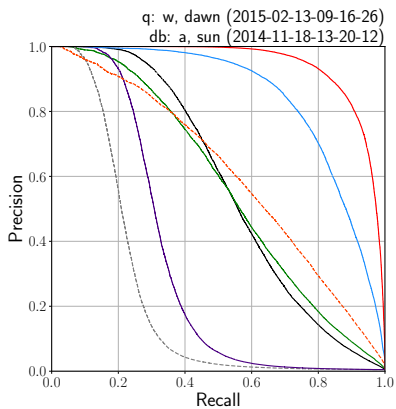
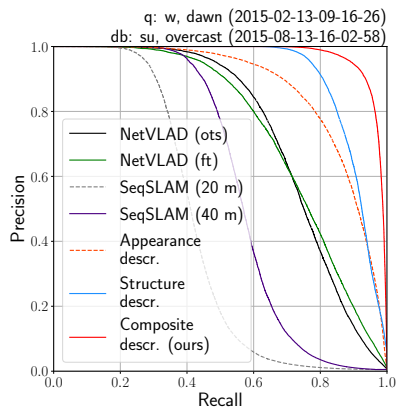
posite descriptors, see Fig. 6. Composite descriptors are produced using the best-performing model architecture which uses $d_A = 12$ and $d_S = 9$ layers for both visual and structural feature extraction networks, respectively, and a simple concatenation of the extracted intermediary features. While descriptors that only encode either appearance or structure often produce incorrect matches, each of the shown queries is correctly matched when using our composite descriptor.

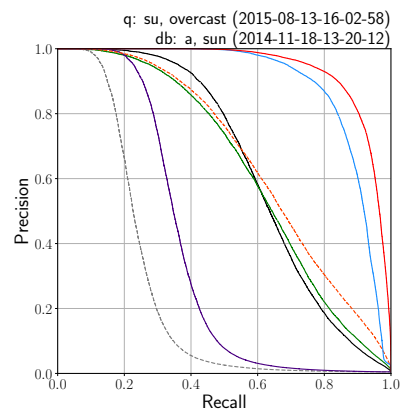
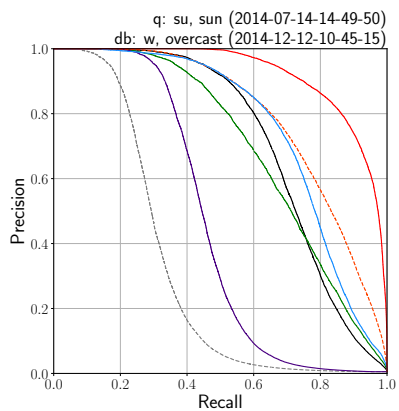
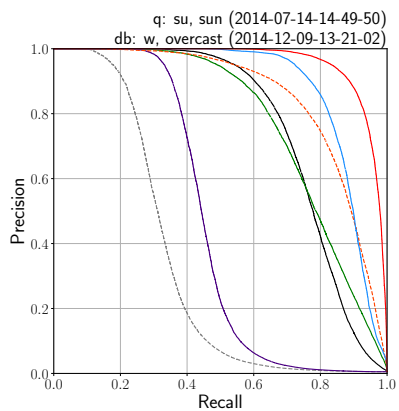
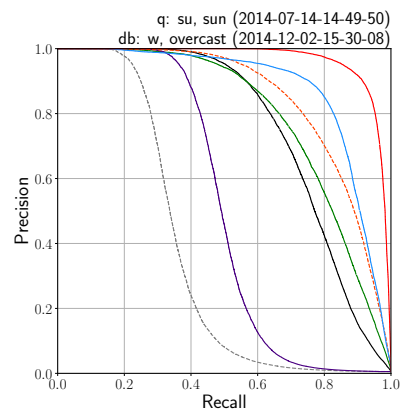
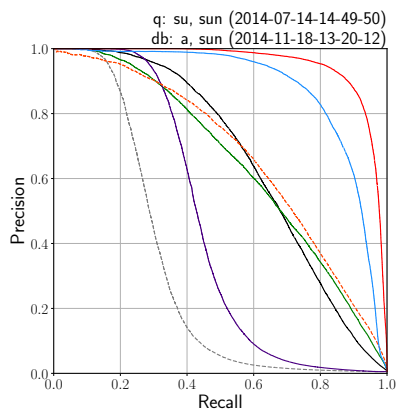
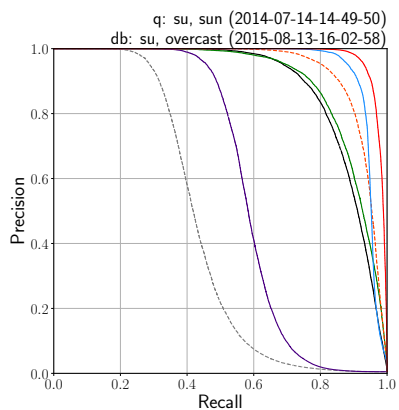
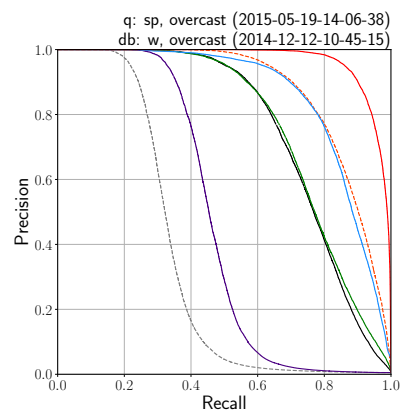
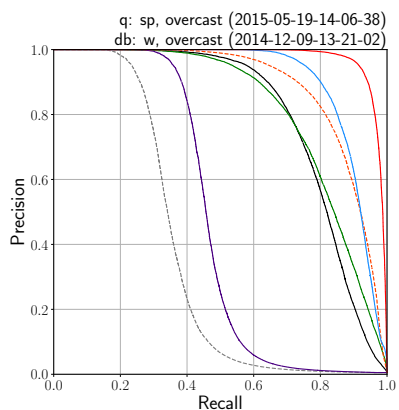
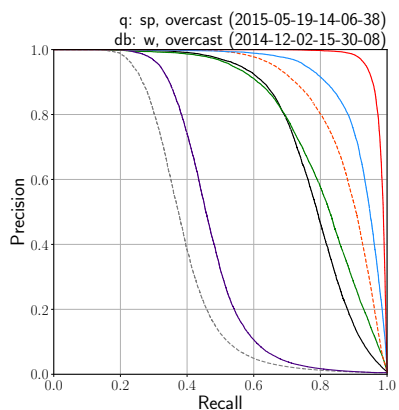
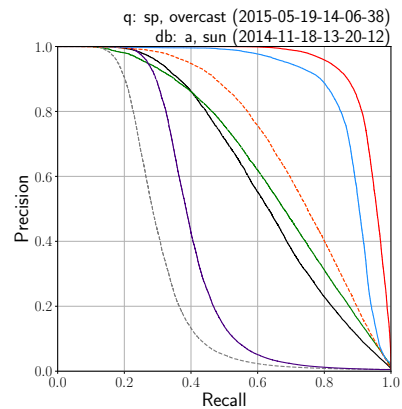
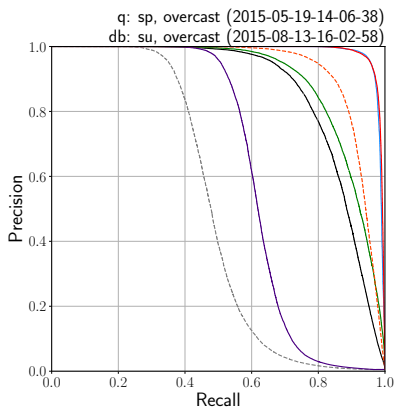
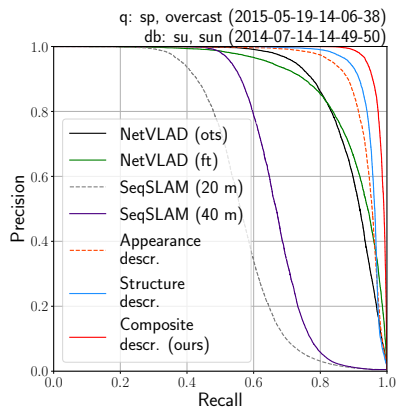
TABLE VII

RECALL FOR $N = 1$ FOR EACH PAIRING OF ROBOTCAR SEQUENCES. A QUERY LOCATION IS DEEMED CORRECTLY MATCHED IF THE RETRIEVED DATABASE INSTANCE EXHIBITS A GROUND TRUTH DISTANCE OF LESS THAN 20 m. EACH SEQUENCE IS DENOTED USING ITS ABBREVIATED TIMESTAMP (*mmddss*). WE UNDERLINE THE BETTER VALUES WHEN COMPARING OFF-THE-SHELF AND FINE-TUNED NETVLAD DESCRIPTORS AS WELL AS APPEARANCE- AND STRUCTURE-BASED DESCRIPTORS. FINE-TUNING NETVLAD IMPROVES ITS RETRIEVAL PERFORMANCE FOR THE MAJORITY OF SEQUENCE PAIRINGS. SUPERIORITY OF USING EITHER VISUAL OR STRUCTURAL CUES VARIES ACROSS SEQUENCE PAIRINGS. HOWEVER, FUSING BOTH MODALITIES INTO OUR COMPOSITE REPRESENTATION RESULTS IN BEST RETRIEVAL PERFORMANCE (PRINTED IN BOLD) ON ALL EXCEPT TWO SEQUENCE PAIRINGS. THE AVERAGES ACROSS ALL SEQUENCE PAIRINGS ARE ALSO REPORTED IN TABLE III.

Query seq.	Database seq.	NetVLAD (ots)	NetVLAD (ft)	MPF	DenseVLAD	SeqSLAM (20 m)	SeqSLAM (40 m)	Appearance descr.	Structural descr.	Composite descr.
020310	021326	0.9669	<u>0.9722</u>	0.9823	0.9678	0.7555	0.8312	0.9751	<u>0.9861</u>	0.9890
020310	031010	<u>0.7355</u>	<u>0.6743</u>	0.7905	0.4845	0.4675	0.5471	<u>0.9780</u>	<u>0.9204</u>	0.9906
020310	051938	0.8584	<u>0.8792</u>	0.9156	0.7902	0.6733	0.7821	<u>0.9306</u>	<u>0.9845</u>	0.9955
020310	071450	<u>0.8438</u>	<u>0.8072</u>	0.7329	0.6287	0.6943	0.7929	0.8744	<u>0.9542</u>	0.9786
020310	081358	<u>0.8473</u>	<u>0.8269</u>	0.9053	0.7478	0.6162	0.6731	<u>0.9302</u>	<u>0.9514</u>	0.9816
020310	111812	<u>0.7518</u>	<u>0.7241</u>	0.7901	0.5743	0.4810	0.6031	<u>0.9086</u>	<u>0.9008</u>	0.9776
020310	120208	0.9359	<u>0.9433</u>	0.9626	0.8857	0.7173	0.7715	<u>0.9792</u>	<u>0.9890</u>	1.0000
020310	120902	0.9249	<u>0.9294</u>	0.9354	0.8849	0.6844	0.7736	0.9755	<u>0.9820</u>	0.9996
020310	121215	0.8429	<u>0.8312</u>	0.8255	0.7886	0.6830	0.7297	0.9161	<u>0.9502</u>	0.9583
021326	031010	<u>0.7984</u>	<u>0.8126</u>	0.8060	0.7293	0.4559	0.5784	<u>0.9792</u>	<u>0.9146</u>	0.9875
021326	051938	0.9159	<u>0.9534</u>	0.9744	0.9584	0.6995	0.8100	<u>0.9525</u>	<u>0.9833</u>	0.9992
021326	071450	0.8948	<u>0.9178</u>	0.7530	0.8864	0.8012	0.8988	0.8645	<u>0.9492</u>	0.9744
021326	081358	0.9159	<u>0.9267</u>	0.9492	0.9279	0.7246	0.8290	0.9559	<u>0.9575</u>	0.9863
021326	111812	0.8247	<u>0.8555</u>	0.8589	0.8442	0.5681	0.6667	<u>0.9413</u>	<u>0.9038</u>	0.9775
021326	120208	0.9721	<u>0.9733</u>	0.9933	0.9842	0.7598	0.8309	<u>0.9988</u>	<u>0.9888</u>	1.0000
021326	120902	<u>0.9867</u>	<u>0.9858</u>	0.9891	0.9913	0.6587	0.7449	<u>0.9896</u>	<u>0.9925</u>	1.0000
021326	121215	0.9012	<u>0.9233</u>	0.8513	0.9007	0.6603	0.6989	0.9122	<u>0.9370</u>	0.9476
031010	051938	0.8063	<u>0.8687</u>	0.8738	0.6218	0.4594	0.5517	0.9219	<u>0.9273</u>	0.9788
031010	071450	0.7808	<u>0.8337</u>	0.6978	0.5452	0.5205	0.6508	<u>0.8887</u>	<u>0.8462</u>	0.9705
031010	081358	0.7984	<u>0.8466</u>	0.8529	0.5777	0.5888	0.7161	<u>0.9273</u>	<u>0.8807</u>	0.9755
031010	111812	0.9252	<u>0.9381</u>	0.9455	0.8433	0.5450	0.6957	<u>0.9472</u>	<u>0.9443</u>	0.9875
031010	120208	0.8246	<u>0.8865</u>	0.9790	0.7510	0.5691	0.6896	<u>0.9850</u>	<u>0.9622</u>	1.0000
031010	120902	0.8263	<u>0.8990</u>	0.9547	0.7469	0.5858	0.6848	<u>0.9618</u>	<u>0.9414</u>	1.0000
031010	121215	0.8075	<u>0.8837</u>	0.8738	0.6839	0.6063	0.7342	<u>0.9297</u>	<u>0.9110</u>	0.9521
051938	071450	0.9579	<u>0.9636</u>	0.7674	0.9602	0.8697	0.8905	0.9411	<u>0.9717</u>	0.9931
051938	081358	0.9674	<u>0.9848</u>	0.9866	0.9770	0.7551	0.7899	0.9904	<u>0.9766</u>	0.9894
051938	111812	0.8898	<u>0.8921</u>	0.9272	0.8696	0.5611	0.6180	<u>0.9284</u>	<u>0.9265</u>	0.9899
051938	120208	0.9316	<u>0.9578</u>	0.9639	0.9490	0.7297	0.7839	<u>0.9816</u>	<u>0.9729</u>	1.0000
051938	120902	0.9500	<u>0.9628</u>	0.9685	0.9477	0.6746	0.7498	<u>0.9587</u>	<u>0.9780</u>	0.9959
051938	121215	0.8662	<u>0.8758</u>	0.8012	0.8156	0.6254	0.6887	0.8657	<u>0.9208</u>	0.9369
071450	081358	0.9721	<u>0.9825</u>	0.9840	0.9787	0.5578	0.6209	0.9628	<u>0.9803</u>	0.9912
071450	111812	0.8774	<u>0.8862</u>	0.9275	0.8101	0.5214	0.6305	0.9502	<u>0.9119</u>	0.9929
071450	120208	0.9321	<u>0.9617</u>	0.9635	0.9371	0.5482	0.5912	<u>0.9332</u>	<u>0.9699</u>	0.9869
071450	120902	0.9398	<u>0.9480</u>	0.9408	0.9130	0.5624	0.6599	0.9294	<u>0.9628</u>	0.9863
071450	121215	0.8363	<u>0.8587</u>	0.7582	0.7194	0.5077	0.5491	0.8065	<u>0.8905</u>	0.9079
081358	111812	0.8842	<u>0.8854</u>	0.9189	0.8306	0.5548	0.6565	0.9293	<u>0.9021</u>	0.9955
081358	120208	0.9290	<u>0.9336</u>	0.9467	0.9311	0.7309	0.8108	<u>0.9633</u>	<u>0.9649</u>	0.9794
081358	120902	0.9479	<u>0.9587</u>	0.9369	0.9326	0.7041	0.7778	0.9517	<u>0.9558</u>	0.9826
081358	121215	0.8504	<u>0.8798</u>	0.7718	0.7618	0.6688	0.7585	0.8669	<u>0.9047</u>	0.9270
111812	120208	0.8854	<u>0.9102</u>	0.9662	0.9043	0.6948	0.8060	0.9580	<u>0.9257</u>	0.9862
111812	120902	0.8764	<u>0.8995</u>	0.9441	0.8638	0.7611	0.8371	<u>0.9336</u>	<u>0.9067</u>	0.9887
111812	121215	0.8018	<u>0.8212</u>	0.8022	0.7367	0.7029	0.8166	<u>0.8748</u>	<u>0.8023</u>	0.9141
120208	120902	0.9901	<u>0.9942</u>	0.9962	0.9909	0.7547	0.8292	<u>0.9954</u>	<u>0.9963</u>	1.0000
120208	121215	0.9251	<u>0.9402</u>	0.8573	0.9156	0.7642	0.8697	0.9051	0.9475	0.9466
120902	121215	0.9316	<u>0.9415</u>	0.8486	0.9225	0.7861	0.8535	0.9023	<u>0.9436</u>	0.9449
Average recall@1		0.8851	0.8999	0.8927	0.8314	0.6447	0.7305	0.9389	0.9421	0.9796







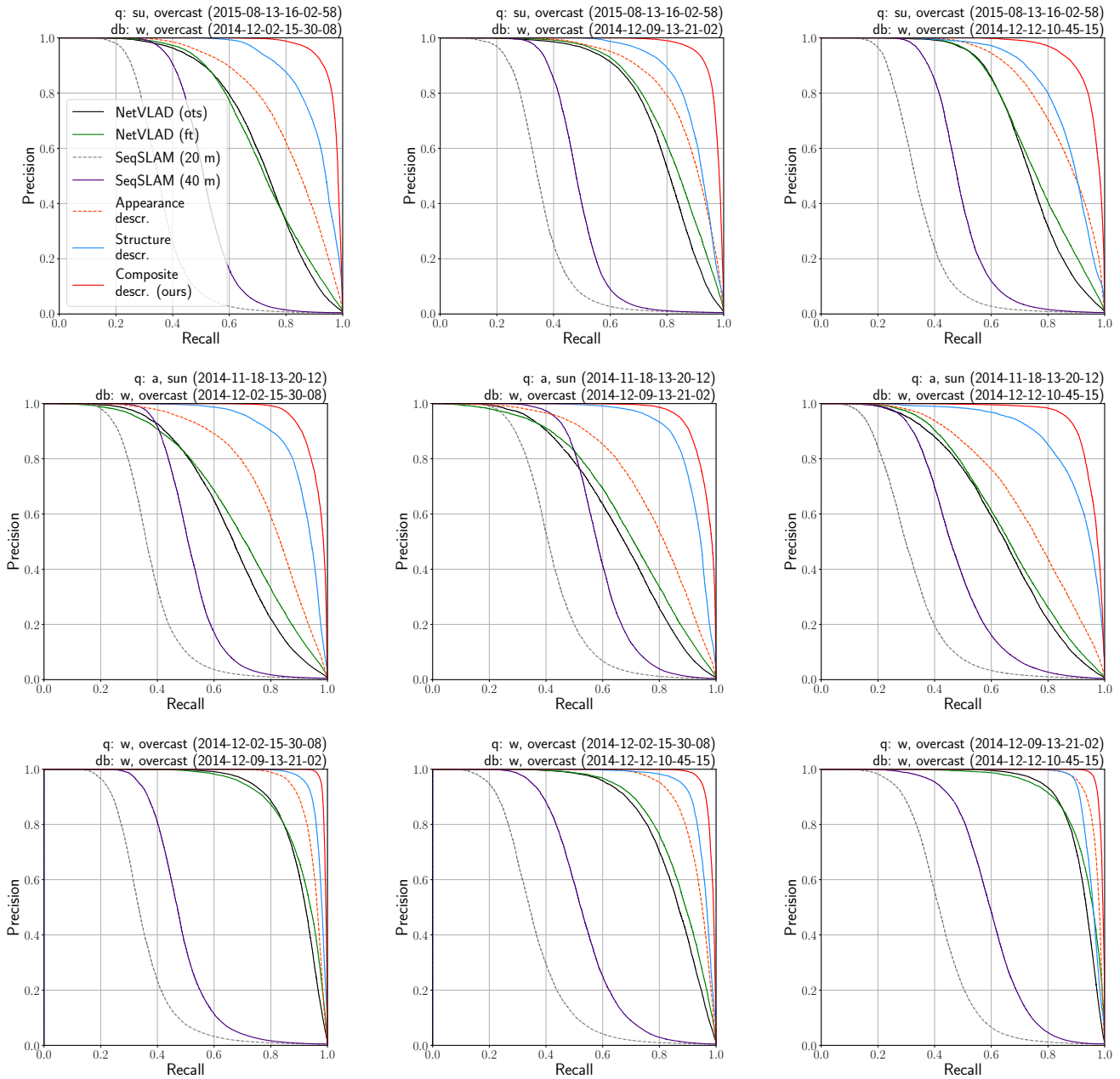


Fig. 5. Precision-recall curves for each pairing of query and database sequences out of a total of 45 sequence combinations. Each curve shows precision and recall, parametrized on the distance threshold $d_{emb,th}$ in embedding space, for exhaustive pairwise matching of query to database sequences. Each plot further details the varying conditions for both query and database sequences which are referenced using their timestamps. We compare our composite descriptors – obtained using a simple concatenation for combining intermediary visual and structural descriptors – to descriptors using only one of the input modalities as well as to off-the-shelf (ots) and fine-tuned (ft) variants of NetVLAD descriptors, as detailed in Section V-B. Specifically, pairwise matching based on our composite descriptor outperforms all other descriptors for *every* sequence combination.



Query

Appearance
descr.

Structure
descr.

Composite
descr. (ours)

Fig. 6. Examples of top retrieved database instances using appearance, structure, or composite descriptors. Only images of the retrieved instances are shown. Green and red borders indicate correct and incorrect matches. While descriptors that separately leverage either appearance or structure often produce incorrect matches, each query is successfully matched using our composite descriptor despite changes in viewpoint and visual appearance.