

Self-Supervised Multi-Object Tracking with Cross-Input Consistency

Favyen Bastani, Songtao He, Sam Madden (MIT CSAIL)



NATIONAL PARK SERVICE
explorer.org

Labeling tracks is expensive



Object Detection
Training Example



Object Tracking
Training Example

Our Approach: Training a Robust Tracker on Unlabeled Video

Object
Detector

+



Robust Multi-Object
Tracking Model

Despite training only on unlabeled video, our unsupervised approach performs competitively with several fully supervised methods that train on video-level labels in MOT17 and Kitti!

Input: Unlabeled Video + Predicted Object Bounding Boxes

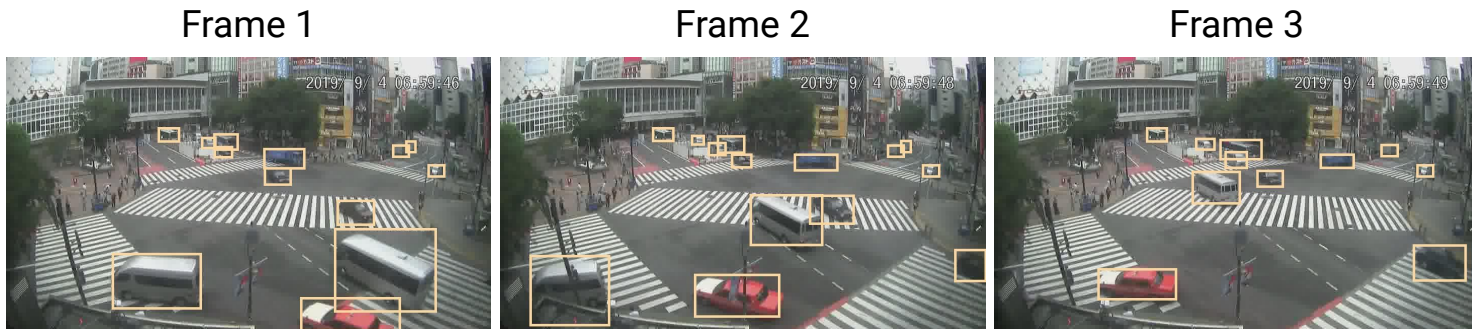


Object
Detector



Input-Hiding Scheme – Visual-Spatial Hiding

Original
Video Sequence



Input 1:
Visual Information



Input 2:
Spatial Information

x=93	x=950	x=652	x=899
y=681	y=459	y=333	y=337
w=176	w=143	w=78	w=62
h=103	h=75	h=87	h=38

x=490	x=740	x=821	x=848
y=429	y=575	y=335	y=325
w=211	w=235	w=63	w=51
h=140	h=127	h=47	h=38

x=458	x=847	x=819	x=609
y=313	y=325	y=336	y=282
w=62	w=53	w=62	w=75
h=49	h=37	h=47	h=48

Input 1:
Visual Information



	(a)	(b)	(c)	(d)
(1)	0.3	0.2	0.0	0.1
(2)	0.0	0.4	0.0	0.2
(3)	0.1	0.0	0.9	0.0
(4)	0.1	0.2	0.1	0.1

Use inconsistency between
matrices as a learning signal!

	(a)	(b)	(c)	(d)
(1)	0.0	0.1	0.2	0.2
(2)	0.1	0.1	0.1	0.1
(3)	0.0	0.3	0.7	0.0
(4)	0.4	0.0	0.1	0.1

Input 2:
Spatial Information

(1)	(2)	(3)	(4)
x=93	x=950	x=652	x=899
y=681	y=459	y=333	y=337
w=176	w=143	w=78	w=62
h=103	h=75	h=87	h=38

x=490	x=740	x=821	x=848
y=429	y=575	y=335	y=325
w=211	w=235	w=63	w=51
h=140	h=127	h=47	h=38

(a)	(b)	(c)	(d)
x=458	x=847	x=819	x=609
y=313	y=325	y=336	y=282
w=62	w=53	w=62	w=75
h=49	h=37	h=47	h=48

Why wouldn't model create arbitrary tracking outputs?

Saturation: 16



Saturation: 18



???

Frame 1

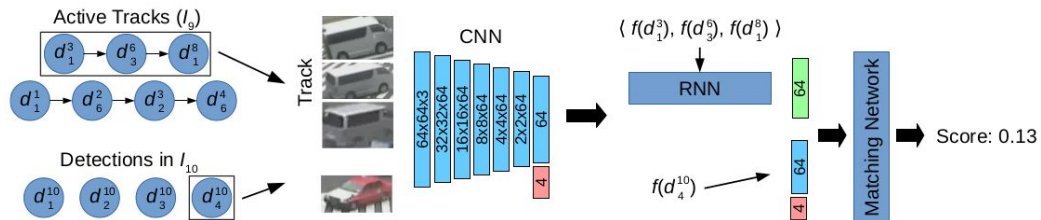
x=93	x=950	x=652	x=899
y=681	y=459	y=333	y=337
w=176	w=143	w=78	w=62
h=103	h=75	h=87	h=38

Frame 2

x=490	x=740	x=821	x=848
y=429	y=575	y=335	y=325
w=211	w=235	w=63	w=51
h=140	h=127	h=47	h=38

Frame 3

x=458	x=847	x=819	x=609
y=313	y=325	y=336	y=282
w=62	w=53	w=62	w=75
h=49	h=37	h=47	h=48



Model Architecture

$$\begin{bmatrix} 0.3 & 0.2 & 0.0 & 0.1 \\ 0.0 & 0.4 & 0.0 & 0.2 \\ 0.1 & 0.0 & 0.9 & 0.0 \\ 0.1 & 0.2 & 0.1 & 0.1 \end{bmatrix}$$

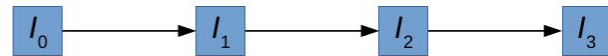
Transition Matrix

```
array([[0.40525553, 0.58487146, 0.35028426],
       [0.17897008, 0.99813136, 0.65189838],
       [0.53625954, 0.51995857, 0.71777389]])
```



```
array([[0.42718929, 0.52618008, 0.84496057],
       [0.8125701, 0.29344281, 0.09825419],
       [0.19191158, 0.11864856, 0.82009976]])
```

**Consistency
Loss Function**



Tracker A (Visual)



Tracker B (Spatial)

(93, 681, 176, 103) (490, 429, 211, 140) (847, 325, 53, 37) (609, 282, 75, 48)
 (950, 459, 143, 75) (740, 575, 235, 127) (458, 313, 62, 49) (899, 337, 62, 48)
 (652, 333, 78, 87) (821, 335, 63, 47) (819, 336, 62, 47) (848, 325, 51, 38)

Input-Hiding Scheme

Frame $i - 2$



Frame $i - 1$



Frame i



Frame $i + 1$



Frame $i + 2$



Frame $i - 2$

Frame $i - 1$

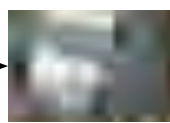
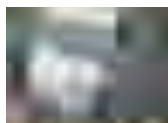
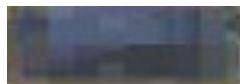
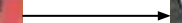
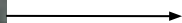
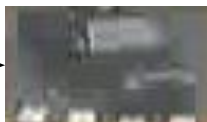
Frame i

Frame $i + 1$

Frame $i + 2$



Track Prefixes



Frame i - 2

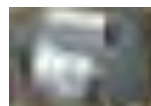
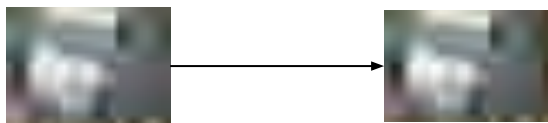
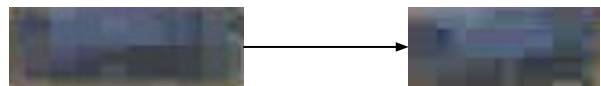
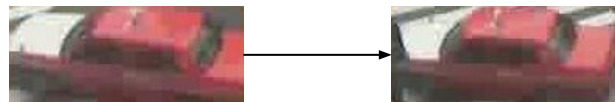
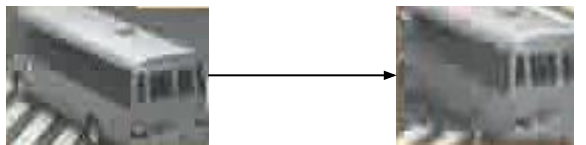
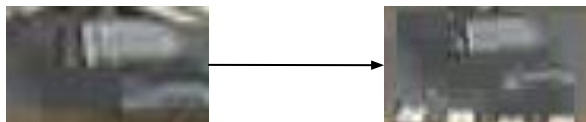
Frame i - 1

Frame i



Active Track Prefixes

Current Detections



Frame $i - 2$

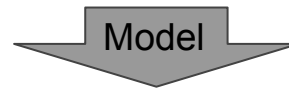
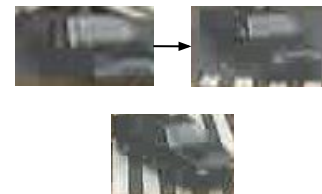
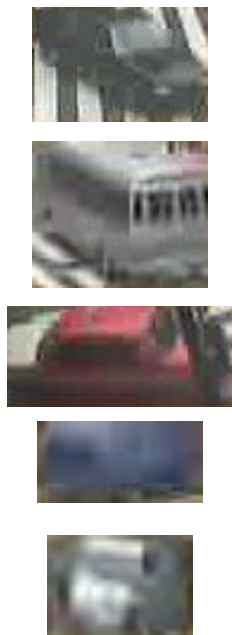
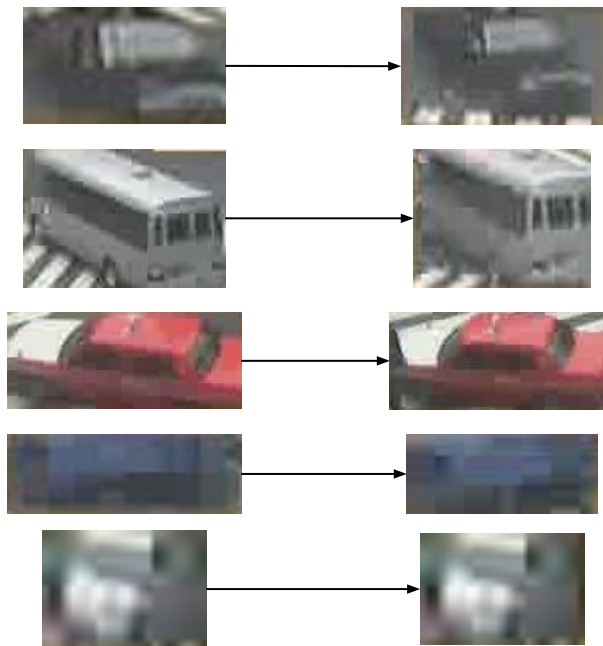
Frame $i - 1$

Frame i

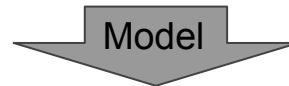


Active Track Prefixes

Current Detections



Same Object



Different Objects

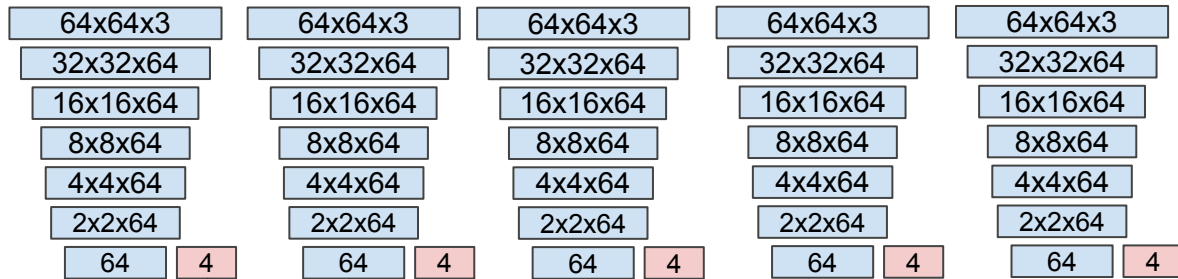


Track Prefix

Current Detection



Track Prefix



Current Detection

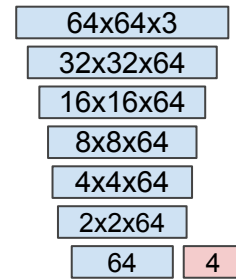
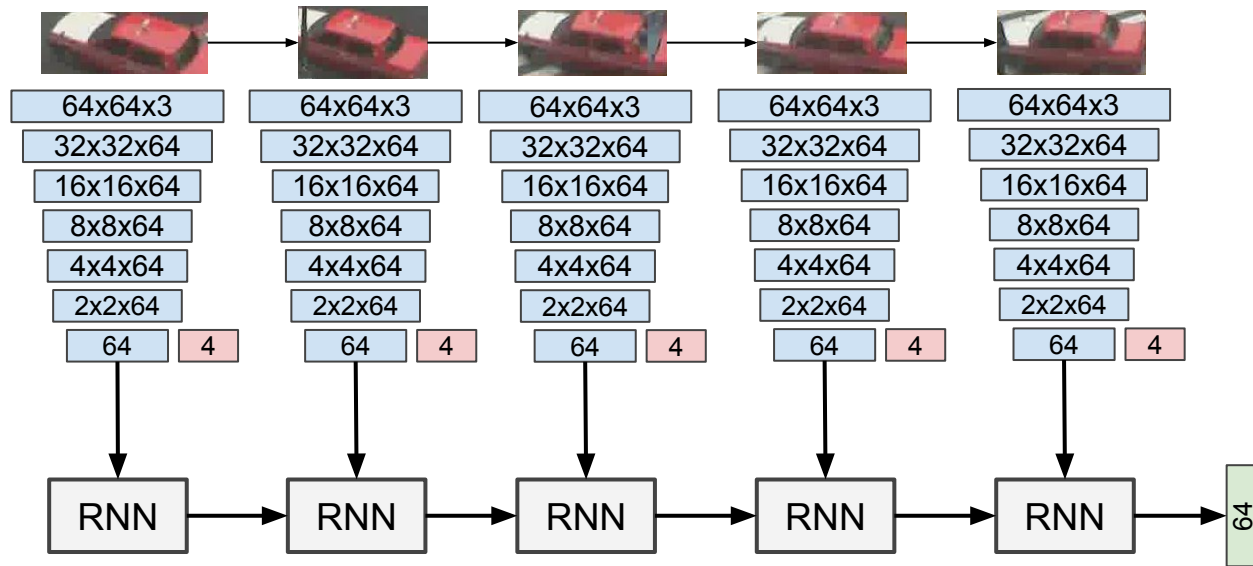


Image Features
from CNN

4D Bounding Box
Coordinates

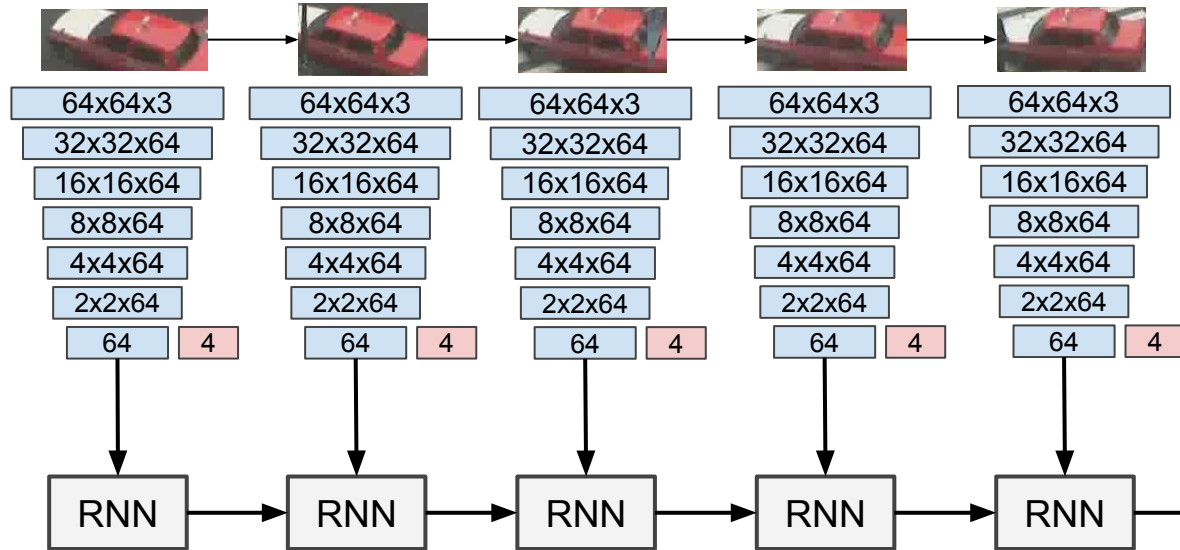
Track Prefix



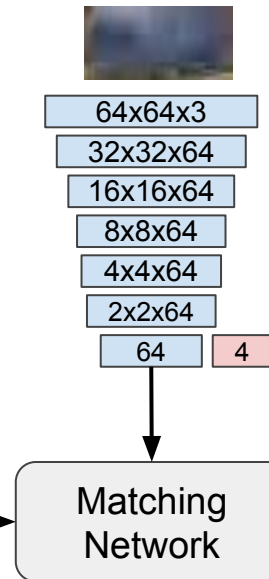
Current Detection



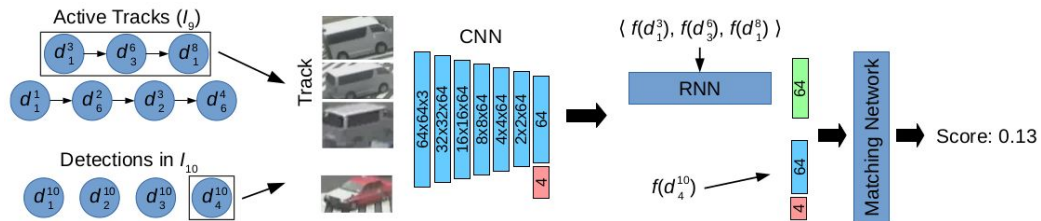
Track Prefix



Current Detection



Score: 0.02
(Probably not the same)



Model Architecture



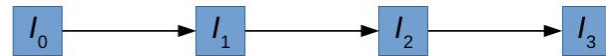
Transition Matrix

```
array([[0.40525553, 0.58487146, 0.35028426],
       [0.17897008, 0.99813136, 0.65189838],
       [0.53625954, 0.51995857, 0.71777389]])
```



```
array([[0.42718929, 0.52618008, 0.84496057],
       [0.8125701, 0.29344281, 0.09825419],
       [0.19191158, 0.11864856, 0.82009976]])
```

Consistency Loss Function



Tracker A (Visual)



Tracker B (Spatial)

(93, 681, 176, 103) (490, 429, 211, 140) (847, 325, 53, 37) (609, 282, 75, 48)
 (950, 459, 143, 75) (740, 575, 235, 127) (458, 313, 62, 49) (899, 337, 62, 48)
 (652, 333, 78, 87) (821, 335, 63, 47) (819, 336, 62, 47) (848, 325, 51, 38)

Input-Hiding Scheme

Training
Example

Frame 1



Frame 2

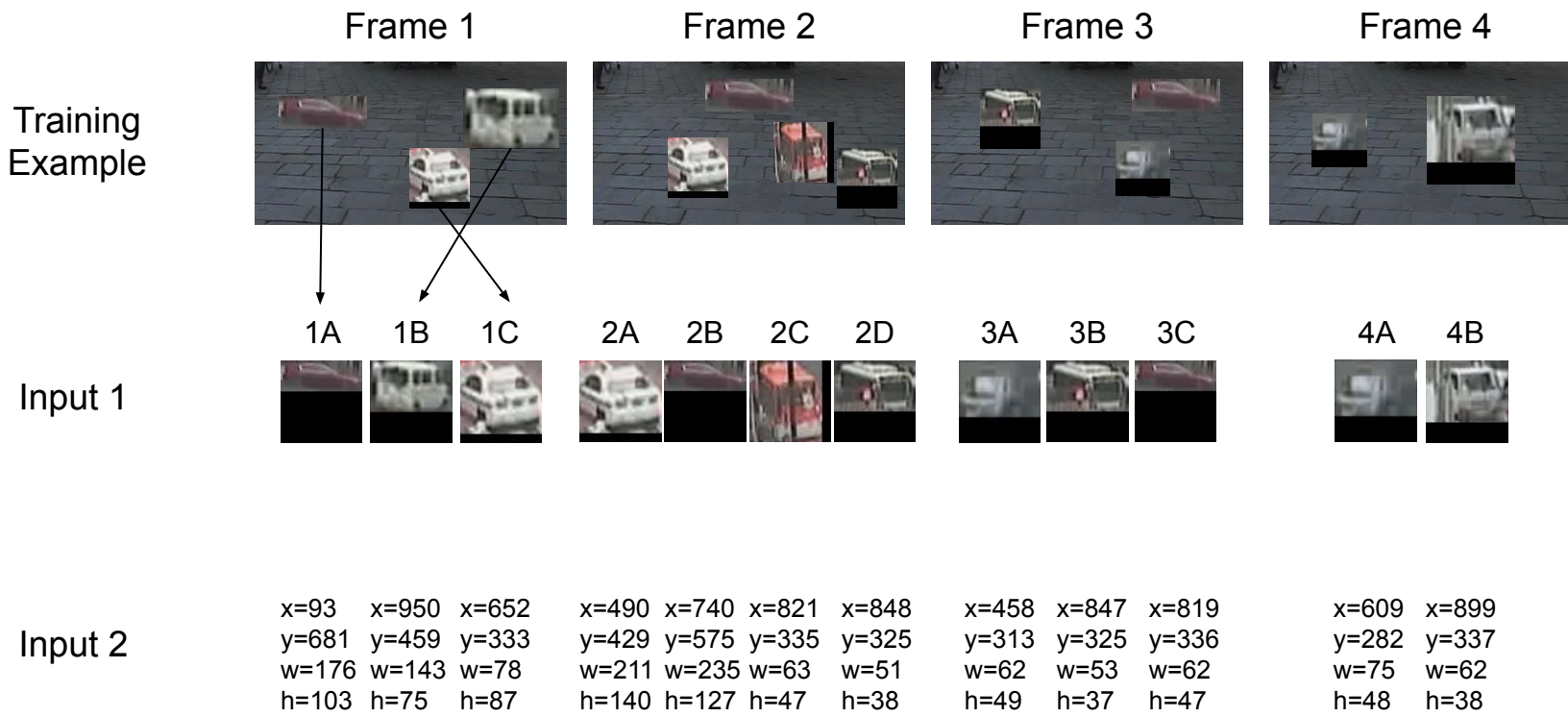


Frame 3

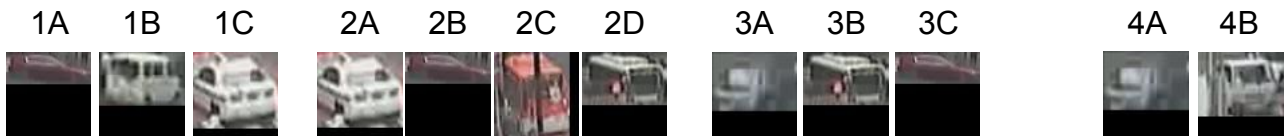


Frame 4





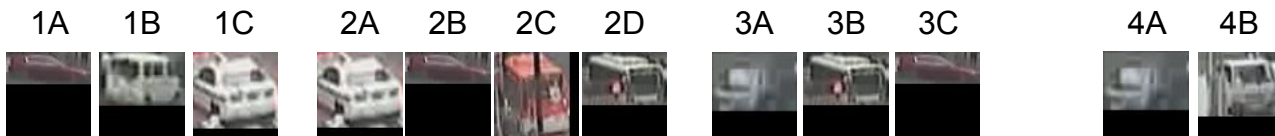
Input 1



Input 2

x=93	x=950	x=652	x=490	x=740	x=821	x=848	x=458	x=847	x=819	x=609	x=899
y=681	y=459	y=333	y=429	y=575	y=335	y=325	y=313	y=325	y=336	y=282	y=337
w=176	w=143	w=78	w=211	w=235	w=63	w=51	w=62	w=53	w=62	w=75	w=62
h=103	h=75	h=87	h=140	h=127	h=47	h=38	h=49	h=37	h=47	h=48	h=38

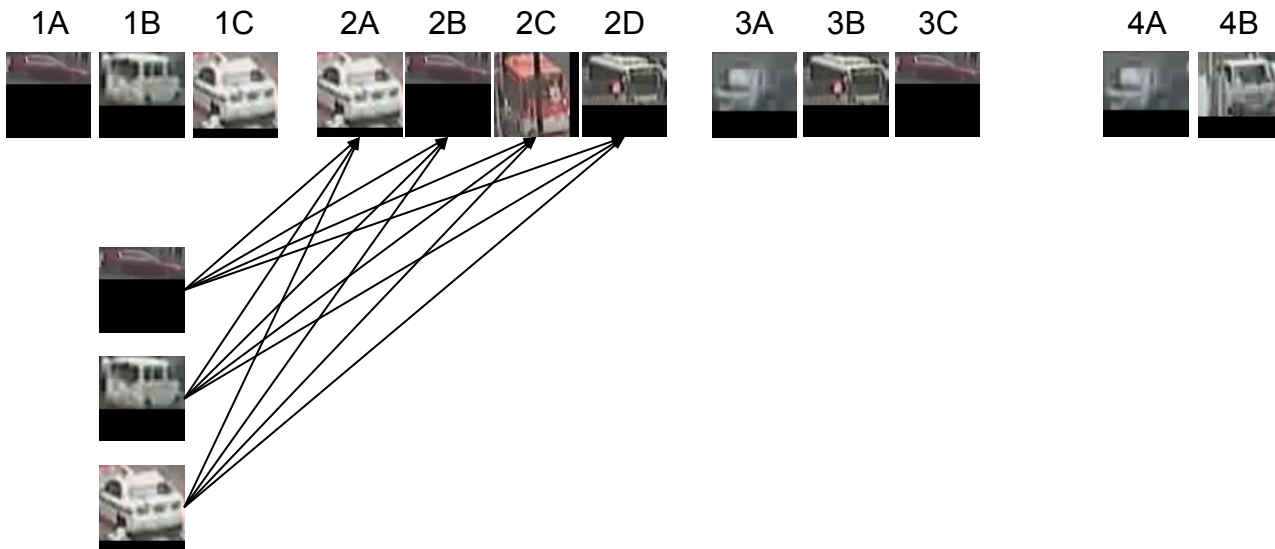
Input 1



Input 2

x=93	x=950	x=652	x=490	x=740	x=821	x=848	x=458	x=847	x=819	x=609	x=899
y=681	y=459	y=333	y=429	y=575	y=335	y=325	y=313	y=325	y=336	y=282	y=337
w=176	w=143	w=78	w=211	w=235	w=63	w=51	w=62	w=53	w=62	w=75	w=62
h=103	h=75	h=87	h=140	h=127	h=47	h=38	h=49	h=37	h=47	h=48	h=38

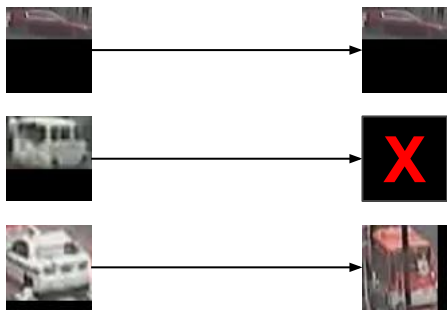
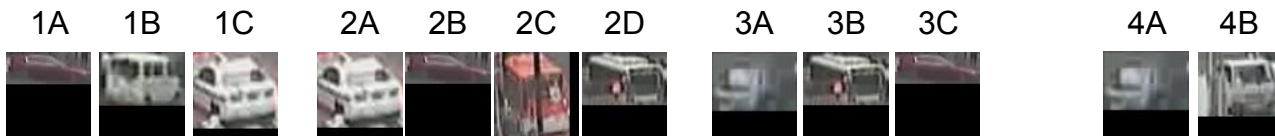
Input 1



Input 2

x=93	x=950	x=652	x=490	x=740	x=821	x=848	x=458	x=847	x=819	x=609	x=899
y=681	y=459	y=333	y=429	y=575	y=335	y=325	y=313	y=325	y=336	y=282	y=337
w=176	w=143	w=78	w=211	w=235	w=63	w=51	w=62	w=53	w=62	w=75	w=62
h=103	h=75	h=87	h=140	h=127	h=47	h=38	h=49	h=37	h=47	h=48	h=38

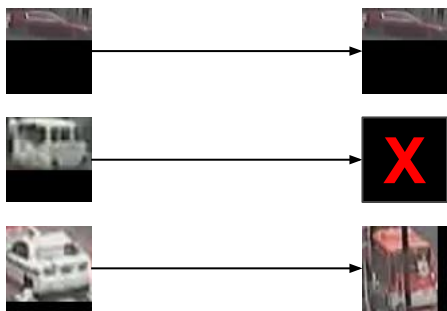
Input 1



Input 2

x=93	x=950	x=652	x=490	x=740	x=821	x=848	x=458	x=847	x=819	x=609	x=899
y=681	y=459	y=333	y=429	y=575	y=335	y=325	y=313	y=325	y=336	y=282	y=337
w=176	w=143	w=78	w=211	w=235	w=63	w=51	w=62	w=53	w=62	w=75	w=62
h=103	h=75	h=87	h=140	h=127	h=47	h=38	h=49	h=37	h=47	h=48	h=38

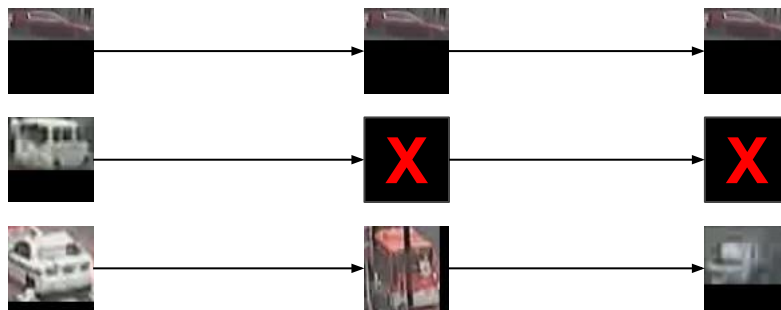
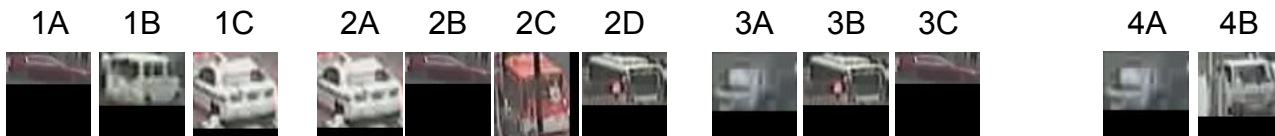
Input 1



Input 2

x=93	x=950	x=652	x=490	x=740	x=821	x=848	x=458	x=847	x=819	x=609	x=899
y=681	y=459	y=333	y=429	y=575	y=335	y=325	y=313	y=325	y=336	y=282	y=337
w=176	w=143	w=78	w=211	w=235	w=63	w=51	w=62	w=53	w=62	w=75	w=62
h=103	h=75	h=87	h=140	h=127	h=47	h=38	h=49	h=37	h=47	h=48	h=38

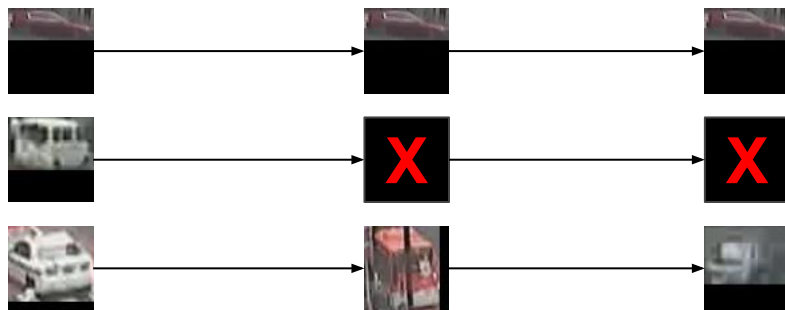
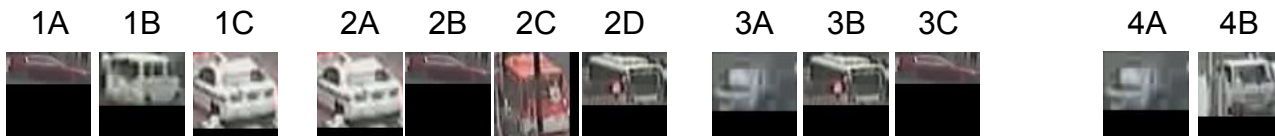
Input 1



Input 2

x=93	x=950	x=652	x=490	x=740	x=821	x=848	x=458	x=847	x=819	x=609	x=899
y=681	y=459	y=333	y=429	y=575	y=335	y=325	y=313	y=325	y=336	y=282	y=337
w=176	w=143	w=78	w=211	w=235	w=63	w=51	w=62	w=53	w=62	w=75	w=62
h=103	h=75	h=87	h=140	h=127	h=47	h=38	h=49	h=37	h=47	h=48	h=38

Input 1

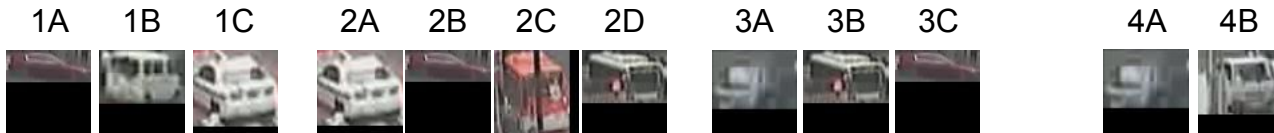


	4A	4B
Track 1A	-23	-2
Track 1B	-1	-1
Track 1C	16	-15

Input 2

x=93	x=950	x=652	x=490	x=740	x=821	x=848	x=458	x=847	x=819	x=609	x=899
y=681	y=459	y=333	y=429	y=575	y=335	y=325	y=313	y=325	y=336	y=282	y=337
w=176	w=143	w=78	w=211	w=235	w=63	w=51	w=62	w=53	w=62	w=75	w=62
h=103	h=75	h=87	h=140	h=127	h=47	h=38	h=49	h=37	h=47	h=48	h=38

Input 1



	4A	4B
Track 1A	-23	-2
Track 1B	-1	-1
Track 1C	16	-15

Input 2

x=93	x=950	x=652	x=490	x=740	x=821	x=848	x=458	x=847	x=819	x=609	x=899
y=681	y=459	y=333	y=429	y=575	y=335	y=325	y=313	y=325	y=336	y=282	y=337
w=176	w=143	w=78	w=211	w=235	w=63	w=51	w=62	w=53	w=62	w=75	w=62
h=103	h=75	h=87	h=140	h=127	h=47	h=38	h=49	h=37	h=47	h=48	h=38



	4A	4B
Track 1A	-23	-2
Track 1B	-1	-1
Track 1C	16	-15

Row-wise Softmax

0.0	1.0
0.5	0.5
1.0	0.0

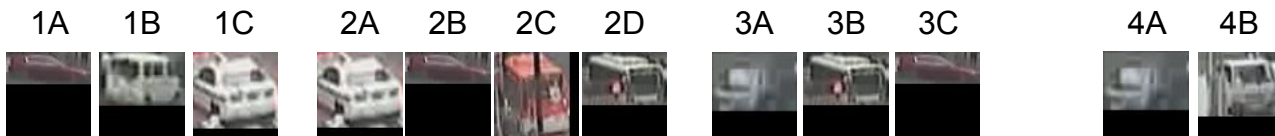
Column-wise Softmax

0.0	0.3
0.0	0.7
1.0	0.0

Input 2

x=93	x=950	x=652	x=490	x=740	x=821	x=848	x=458	x=847	x=819	x=609	x=899
y=681	y=459	y=333	y=429	y=575	y=335	y=325	y=313	y=325	y=336	y=282	y=337
w=176	w=143	w=78	w=211	w=235	w=63	w=51	w=62	w=53	w=62	w=75	w=62
h=103	h=75	h=87	h=140	h=127	h=47	h=38	h=49	h=37	h=47	h=48	h=38

Input 1



	4A	4B
Track 1A	-23	-2
Track 1B	-1	-1
Track 1C	16	-15

Row-wise Softmax

$$\begin{bmatrix} 0.0 & 1.0 \\ 0.5 & 0.5 \\ 1.0 & 0.0 \end{bmatrix}$$

Column-wise Softmax

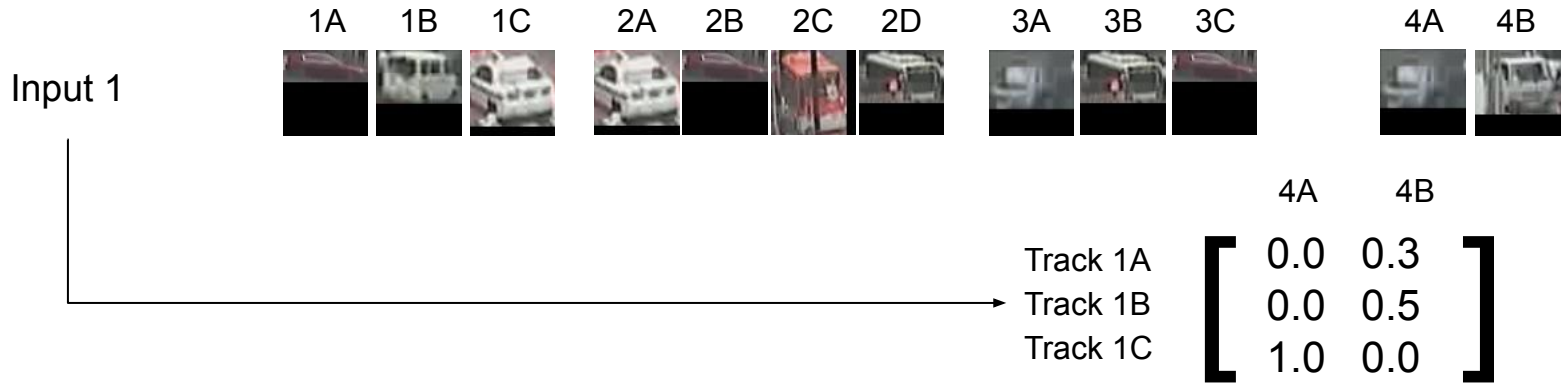
$$\begin{bmatrix} 0.0 & 0.3 \\ 0.0 & 0.7 \\ 1.0 & 0.0 \end{bmatrix}$$

↓ ↓

$$\begin{bmatrix} 0.0 & 0.3 \\ 0.0 & 0.5 \\ 1.0 & 0.0 \end{bmatrix}$$

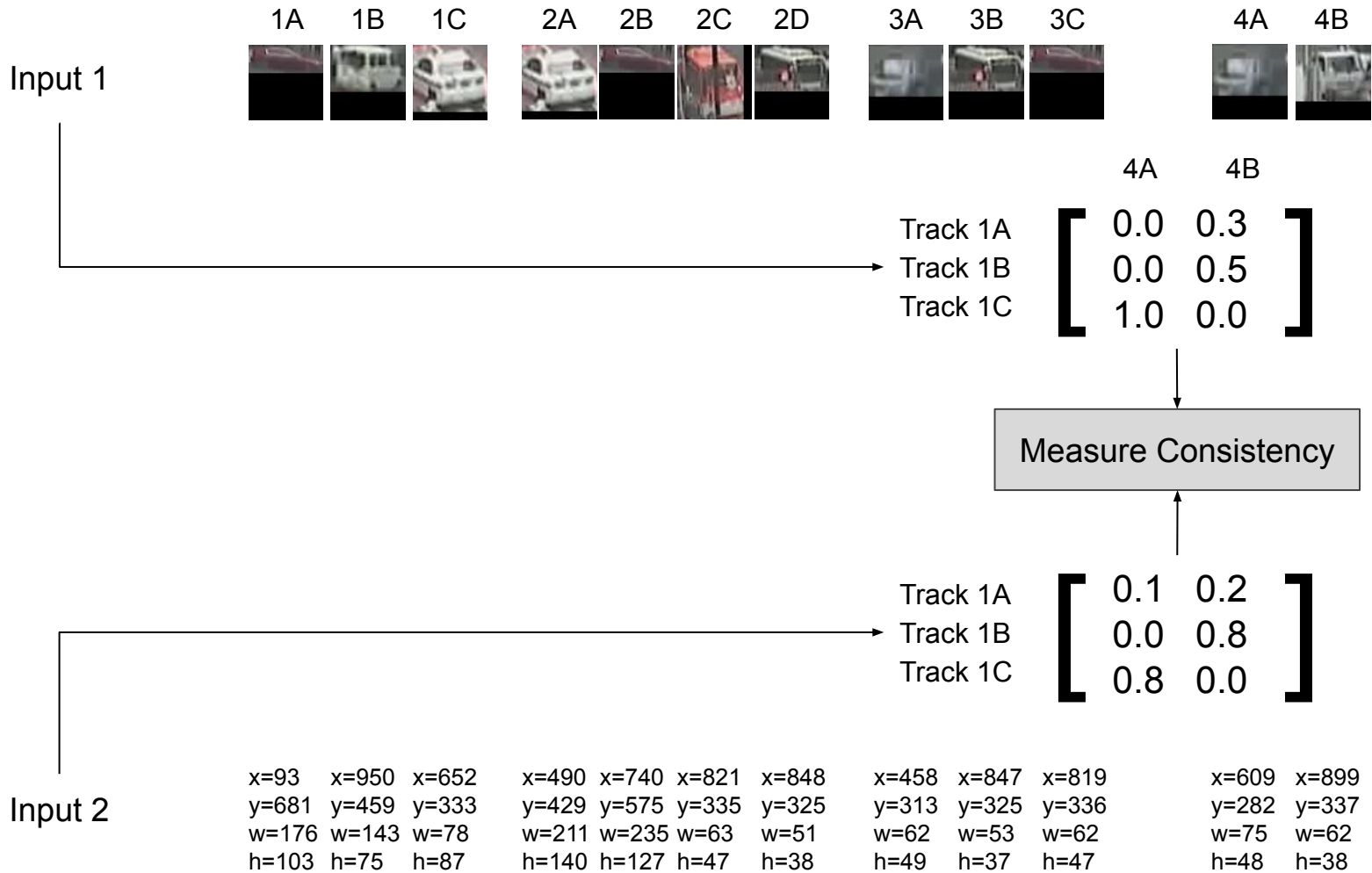
Input 2

x=93	x=950	x=652	x=490	x=740	x=821	x=848	x=458	x=847	x=819	x=609	x=899
y=681	y=459	y=333	y=429	y=575	y=335	y=325	y=313	y=325	y=336	y=282	y=337
w=176	w=143	w=78	w=211	w=235	w=63	w=51	w=62	w=53	w=62	w=75	w=62
h=103	h=75	h=87	h=140	h=127	h=47	h=38	h=49	h=37	h=47	h=48	h=38

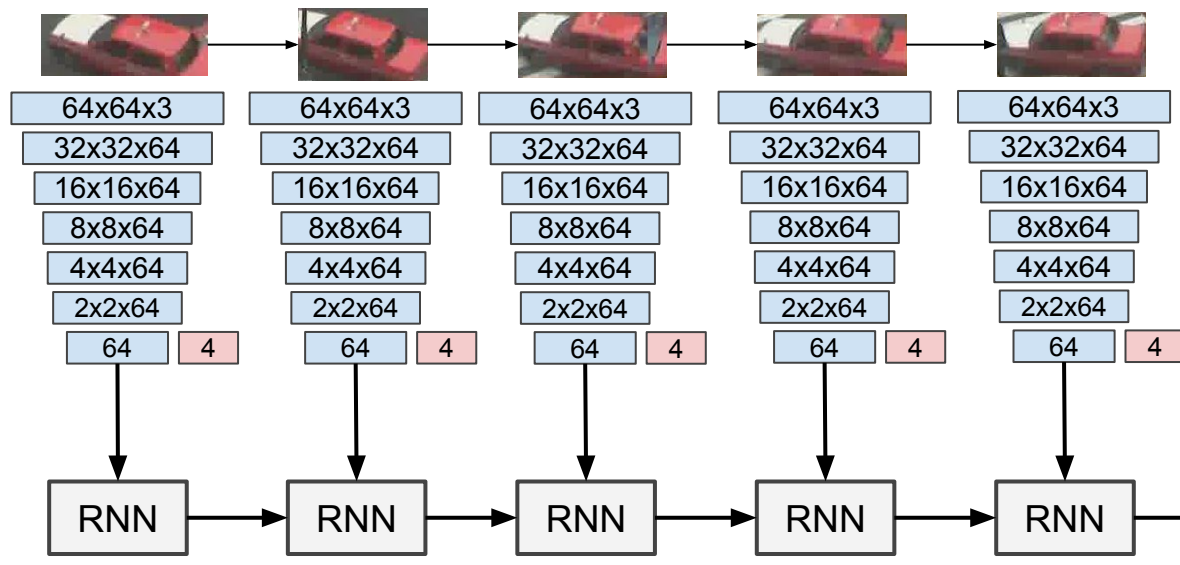


Input 2

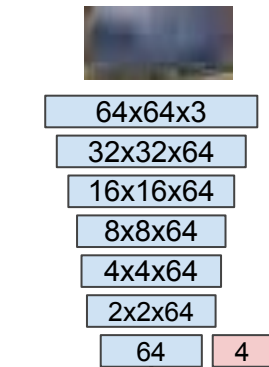
x=93	x=950	x=652	x=490	x=740	x=821	x=848	x=458	x=847	x=819	x=609	x=899
y=681	y=459	y=333	y=429	y=575	y=335	y=325	y=313	y=325	y=336	y=282	y=337
w=176	w=143	w=78	w=211	w=235	w=63	w=51	w=62	w=53	w=62	w=75	w=62
h=103	h=75	h=87	h=140	h=127	h=47	h=38	h=49	h=37	h=47	h=48	h=38



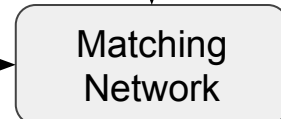
Track Prefix



Current Detection

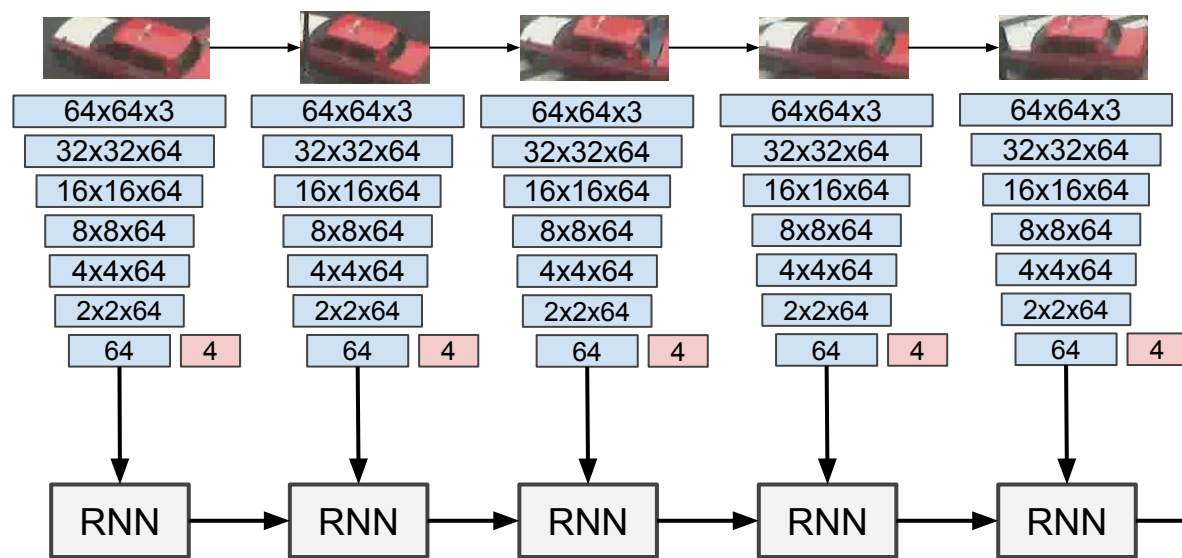


64

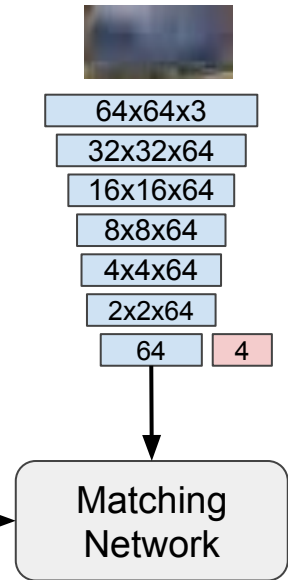


Score: 0.02
(Probably not the same)

Track Prefix



Current Detection



Score: 0.02
(Probably not the same)

4A 4B

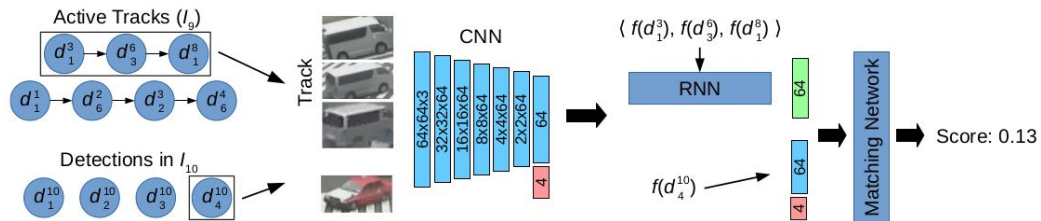
Track 1A	0.0	0.3
Track 1B	0.0	0.5
Track 1C	1.0	0.0

Input 1

4A 4B

Track 1A	0.0	0.3
Track 1B	0.0	0.5
Track 1C	1.0	0.0

Input 2



Model Architecture

$$\begin{bmatrix} 0.3 & 0.2 & 0.0 & 0.1 \\ 0.0 & 0.4 & 0.0 & 0.2 \\ 0.1 & 0.0 & 0.9 & 0.0 \\ 0.1 & 0.2 & 0.1 & 0.1 \end{bmatrix}$$

Transition Matrix

```
array([[0.40525553, 0.58487146, 0.35028426],
       [0.17897008, 0.99813136, 0.65189838],
       [0.53625954, 0.51995857, 0.71777389]])
```



```
array([[0.42718929, 0.52618008, 0.84496057],
       [0.8125701, 0.29344281, 0.09825419],
       [0.19191158, 0.11864856, 0.82009976]])
```

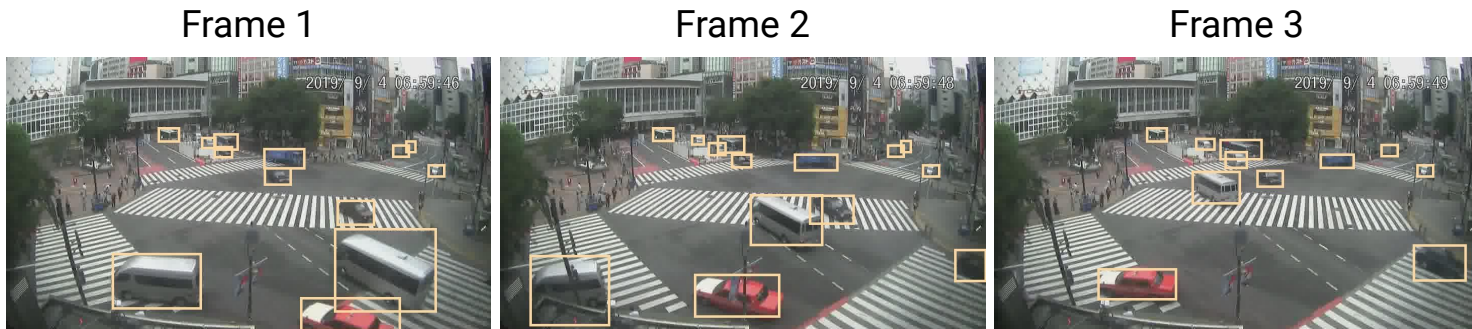
**Consistency
Loss Function**



Input-Hiding Scheme

Input-Hiding Scheme – Visual-Spatial Hiding

Original
Video Sequence



Input 1:
Visual Information

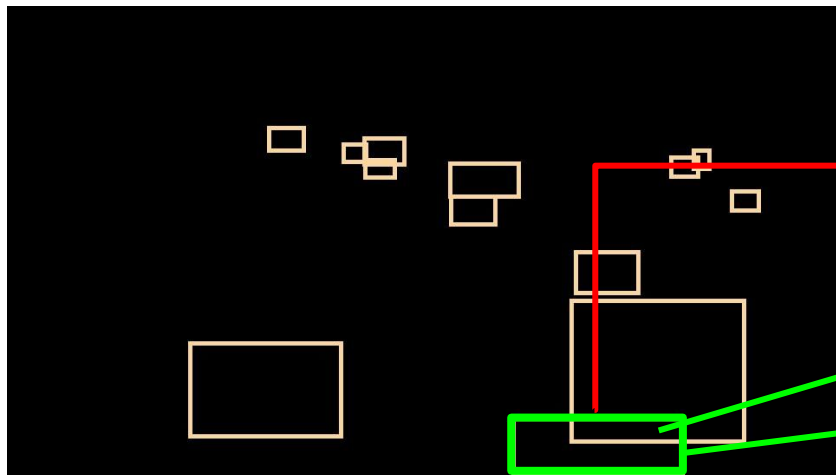


Input 2:
Spatial Information

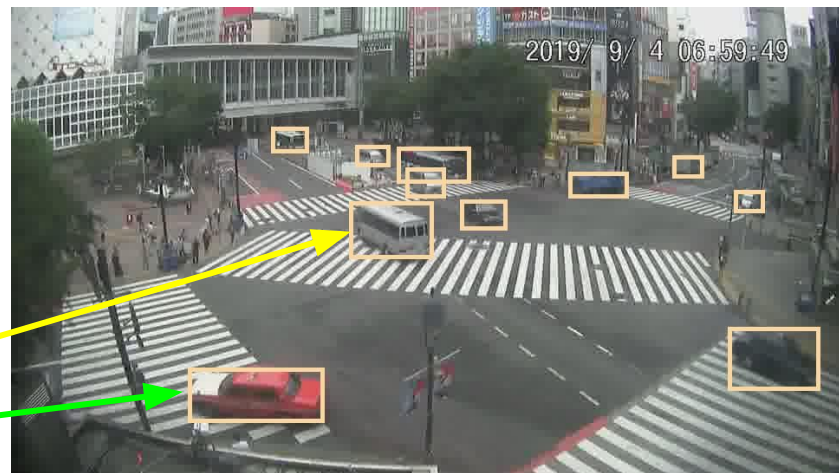
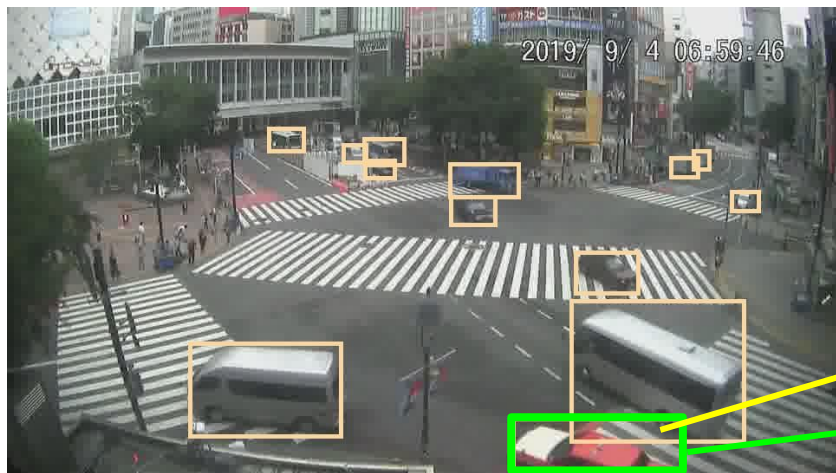
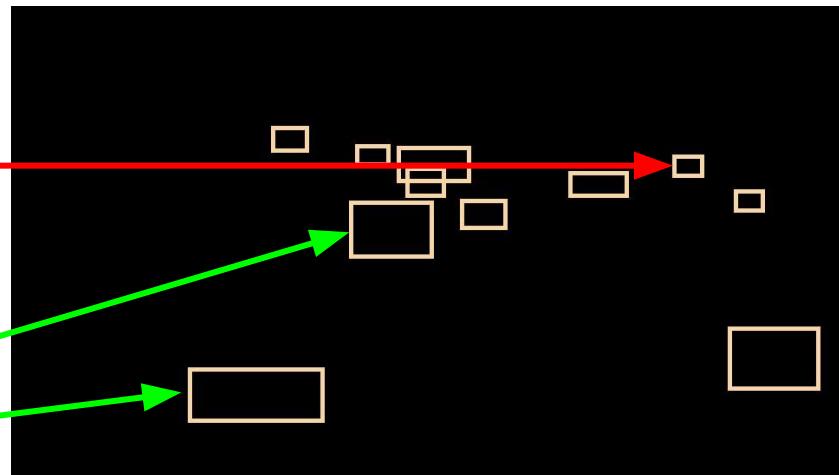
x=93	x=950	x=652	x=899
y=681	y=459	y=333	y=337
w=176	w=143	w=78	w=62
h=103	h=75	h=87	h=38

x=490	x=740	x=821	x=848
y=429	y=575	y=335	y=325
w=211	w=235	w=63	w=51
h=140	h=127	h=47	h=38

x=458	x=847	x=819	x=609
y=313	y=325	y=336	y=282
w=62	w=53	w=62	w=75
h=49	h=37	h=47	h=48



X



Results on MOT17 Benchmark

	Method	IDF1	MOTA	MT	ML	FP	FN	Idsw	Frag
Unsupervised Methods	Visual-Spatial (ours)	58.3	56.8	538	880	12K	231K	1K	2K
	SORT [2]	39.8	43.1	295	997	28K	288K	5K	7K
	IOU [3]	39.4	45.5	369	953	20K	282K	6K	7K
Supervised Methods	Tracktor++ [1]	52.3	53.5	459	861	12K	248K	2K	5K
	MHT-BLSTM [12]	51.9	47.5	429	981	26K	268K	2K	3K
	FAMNet [5]	48.7	52.0	450	787	14K	254K	3K	5K
	LSST [8]	62.3	54.7	480	944	26K	228K	1K	4K
	GSM [18]	57.8	56.4	523	813	14K	230K	1K	3K
	CenterTrack [30]	59.6	61.5	621	752	14K	201K	3K	5K

Table 2: Performance on the MOT17 test set. We compare methods in terms of IDF1 and MOTA, but include other non-comprehensive metrics from MOT17 as well for completeness.

Results on MOT17 Benchmark

	Method	IDF1	MOTA	MT	ML	FP	FN	Idsw	Frag
Unsupervised Methods	Visual-Spatial (ours)	58.3	56.8	538	880	12K	231K	1K	2K
	SORT [2]	39.8	43.1	295	997	28K	288K	5K	7K
	IOU [3]	39.4	45.5	369	953	20K	282K	6K	7K
Supervised Methods	Tracktor++ [1]	52.3	53.5	459	861	12K	248K	2K	5K
	MHT-BLSTM [12]	51.9	47.5	429	981	26K	268K	2K	3K
	FAMNet [5]	48.7	52.0	450	787	14K	254K	3K	5K
	LSST [8]	62.3	54.7	480	944	26K	228K	1K	4K
	GSM [18]	57.8	56.4	523	813	14K	230K	1K	3K
	CenterTrack [30]	59.6	61.5	621	752	14K	201K	3K	5K

Table 2: Performance on the MOT17 test set. We compare methods in terms of IDF1 and MOTA, but include other non-comprehensive metrics from MOT17 as well for completeness.

Results on MOT17 Benchmark

	Method	IDF1	MOTA	MT	ML	FP	FN	Idsw	Frag
Unsupervised Methods	Visual-Spatial (ours)	58.3	56.8	538	880	12K	231K	1K	2K
	SORT [2]	39.8	43.1	295	997	28K	288K	5K	7K
	IOU [3]	39.4	45.5	369	953	20K	282K	6K	7K
Supervised Methods	Tracktor++ [1]	52.3	53.5	459	861	12K	248K	2K	5K
	MHT-BLSTM [12]	51.9	47.5	429	981	26K	268K	2K	3K
	FAMNet [5]	48.7	52.0	450	787	14K	254K	3K	5K
	LSST [8]	62.3	54.7	480	944	26K	228K	1K	4K
	GSM [18]	57.8	56.4	523	813	14K	230K	1K	3K
	CenterTrack [30]	59.6	61.5	621	752	14K	201K	3K	5K

Table 2: Performance on the MOT17 test set. We compare methods in terms of IDF1 and MOTA, but include other non-comprehensive metrics from MOT17 as well for completeness.

Conclusion

- Labeling data for MOT is tedious and costly
- Huge amounts of unlabeled video are almost always available!
- By leveraging unlabeled video, cross-input consistency ***outperforms several fully-supervised MOT methods*** from within the last 1-2 years
- For code and more info: <https://favyen.com/uns20/>

