UUCT - HyMP: Towards Tracking Dispersed Crowd Groups from UAVs

Tonmoay Deb, Mahieyin Rahmun, Shahriar Ali Bijoy, Mayamin Hamid Raha, and Dr. Mohammad Ashrafuzzaman Khan

Department of Electrical & Computer Engineering, North South University

International Joint Conference on Neural Networks 2021 Paper ID: 1827



э

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト







Introduction and Motivation

Simulator

UUCT - dataset generation



æ

▲□▶ ▲圖▶ ▲圖▶ ▲圖▶

Introduction and Motivation

Simulator

- UUCT dataset generation
- Ground truth generation and evaluation



æ

ヘロト ヘロト ヘヨト ヘヨト

Introduction and Motivation

Simulator

- UUCT dataset generation
- Ground truth generation and evaluation
- Proposed tracking algorithm: HyMP



э

Introduction and Motivation

- Simulator
 - UUCT dataset generation
 - Ground truth generation and evaluation
- Proposed tracking algorithm: HyMP
- Results



э

Introduction and Motivation

- Simulator
 - UUCT dataset generation
 - Ground truth generation and evaluation
- Proposed tracking algorithm: HyMP
- Results
- Conclusion
- Simulator Demo



э

Introduction and Motivation

- Simulator
 - UUCT dataset generation
 - Ground truth generation and evaluation
- Proposed tracking algorithm: HyMP
- Results
- Conclusion
- Simulator Demo



э

Introduction and Motivation



A challenging problem: how to track groups of dispersed crowds from aerial viewpoint?



- A challenging problem: how to track groups of dispersed crowds from aerial viewpoint?
- Existing Object Trackers may not be able to adapt to our problem.



- A challenging problem: how to track groups of dispersed crowds from aerial viewpoint?
- Existing Object Trackers may not be able to adapt to our problem.
- Issues with Single Object Trackers (SOT)

- A challenging problem: how to track groups of dispersed crowds from aerial viewpoint?
- Existing Object Trackers may not be able to adapt to our problem.
- Issues with Single Object Trackers (SOT)
 - It can not adapt to crowd group deformations or reformations because they might move sparse.



- A challenging problem: how to track groups of dispersed crowds from aerial viewpoint?
- Existing Object Trackers may not be able to adapt to our problem.
- Issues with Single Object Trackers (SOT)
 - It can not adapt to crowd group deformations or reformations because they might move sparse.
 - The tracker memory can not scale according to the crowd shape changes



- A challenging problem: how to track groups of dispersed crowds from aerial viewpoint?
- Existing Object Trackers may not be able to adapt to our problem.
- Issues with Single Object Trackers (SOT)
 - It can not adapt to crowd group deformations or reformations because they might move sparse.
 - The tracker memory can not scale according to the crowd shape changes
- Issues with Multi-Object Trackers (MOT)



э

- A challenging problem: how to track groups of dispersed crowds from aerial viewpoint?
- Existing Object Trackers may not be able to adapt to our problem.
- Issues with Single Object Trackers (SOT)
 - It can not adapt to crowd group deformations or reformations because they might move sparse.
 - The tracker memory can not scale according to the crowd shape changes
- Issues with Multi-Object Trackers (MOT)
 - Our dataset will have a number of people to track. MOTs will fail to scale tracking each person simultaneously.



э

- A challenging problem: how to track groups of dispersed crowds from aerial viewpoint?
- Existing Object Trackers may not be able to adapt to our problem.
- Issues with Single Object Trackers (SOT)
 - It can not adapt to crowd group deformations or reformations because they might move sparse.
 - The tracker memory can not scale according to the crowd shape changes
- Issues with Multi-Object Trackers (MOT)
 - Our dataset will have a number of people to track. MOTs will fail to scale tracking each person simultaneously.
 - It may not be able to handle crowd dispersion, i.e., split into several groups suddenly because it can not connect the relation between 'multiple targets.'



э

・ロト ・ 『 ト ・ ヨ ト ・ ヨ ト

- A challenging problem: how to track groups of dispersed crowds from aerial viewpoint?
- Existing Object Trackers may not be able to adapt to our problem.
- Issues with Single Object Trackers (SOT)
 - It can not adapt to crowd group deformations or reformations because they might move sparse.
 - The tracker memory can not scale according to the crowd shape changes
- Issues with Multi-Object Trackers (MOT)
 - Our dataset will have a number of people to track. MOTs will fail to scale tracking each person simultaneously.
 - It may not be able to handle crowd dispersion, i.e., split into several groups suddenly because it can not connect the relation between 'multiple targets.'

To tackle these issues, we propose a novel dataset and benchmark, followed by a tracking algorithm that learns end-to-end.



Dataset



<ロト <回ト < 注ト < 注ト

Dataset

Existing datasets deal with object tracking in a generalized manner e.g. not crowd specific



Dataset

- Existing datasets deal with object tracking in a generalized manner e.g. not crowd specific
- Lack of crowd group tracking benchmarks from aerial viewpoint



э

ヘロト ヘロト ヘビト ヘビト

Dataset

- Existing datasets deal with object tracking in a generalized manner e.g. not crowd specific
- Lack of crowd group tracking benchmarks from aerial viewpoint
- The proposed dataset and benchmark: UUCT



э

Dataset

- Existing datasets deal with object tracking in a generalized manner e.g. not crowd specific
- Lack of crowd group tracking benchmarks from aerial viewpoint
- The proposed dataset and benchmark: UUCT
 - Unreal UAV Crowd Tracking



э

Dataset

- Existing datasets deal with object tracking in a generalized manner e.g. not crowd specific
- Lack of crowd group tracking benchmarks from aerial viewpoint
- The proposed dataset and benchmark: UUCT
 - Unreal UAV Crowd Tracking
 - Built with Unreal Engine

э

Dataset

- Existing datasets deal with object tracking in a generalized manner e.g. not crowd specific
- Lack of crowd group tracking benchmarks from aerial viewpoint
- The proposed dataset and benchmark: UUCT
 - Unreal UAV Crowd Tracking
 - Built with Unreal Engine
 - Baseline benchmark



э

Dataset

- Existing datasets deal with object tracking in a generalized manner e.g. not crowd specific
- Lack of crowd group tracking benchmarks from aerial viewpoint
- The proposed dataset and benchmark: UUCT
 - Unreal UAV Crowd Tracking
 - Built with Unreal Engine
 - Baseline benchmark
 - Photorealistic synthetic images



э

Dataset

- Existing datasets deal with object tracking in a generalized manner e.g. not crowd specific
- Lack of crowd group tracking benchmarks from aerial viewpoint
- The proposed dataset and benchmark: UUCT
 - Unreal UAV Crowd Tracking
 - Built with Unreal Engine
 - Baseline benchmark
 - Photorealistic synthetic images
 - Automated ground truth generation



Dataset

- Existing datasets deal with object tracking in a generalized manner e.g. not crowd specific
- Lack of crowd group tracking benchmarks from aerial viewpoint
- The proposed dataset and benchmark: UUCT
 - Unreal UAV Crowd Tracking
 - Built with Unreal Engine
 - Baseline benchmark
 - Photorealistic synthetic images
 - Automated ground truth generation
 - Rich in crowd-specific attributes



Dataset

- Existing datasets deal with object tracking in a generalized manner e.g. not crowd specific
- Lack of crowd group tracking benchmarks from aerial viewpoint
- The proposed dataset and benchmark: UUCT
 - Unreal UAV Crowd Tracking
 - Built with Unreal Engine
 - Baseline benchmark
 - Photorealistic synthetic images
 - Automated ground truth generation
 - Rich in crowd-specific attributes
 - Evaluated on existing tracking algorithms and HyMP



э

Simulator



Built with Unreal Engine 4



<ロト <回ト < 回ト < 回ト

- Built with Unreal Engine 4
- Uses Microsoft AirSim plugin for UAV deployment, control and image capturing



- Built with Unreal Engine 4
- Uses Microsoft AirSim plugin for UAV deployment, control and image capturing
- Al controlled NPCs for directed/randomized movement around the levels



э

ヘロト ヘロト ヘヨト ヘヨト

- Built with Unreal Engine 4
- Uses Microsoft AirSim plugin for UAV deployment, control and image capturing
- AI controlled NPCs for directed/randomized movement around the levels
- A baseline benchmark



э

ヘロト ヘロト ヘビト ヘビト

- Built with Unreal Engine 4
- Uses Microsoft AirSim plugin for UAV deployment, control and image capturing
- AI controlled NPCs for directed/randomized movement around the levels
- A baseline benchmark
 - No environmental clutters/props



э

A D > A P > A D > A D >

- Built with Unreal Engine 4
- Uses Microsoft AirSim plugin for UAV deployment, control and image capturing
- Al controlled NPCs for directed/randomized movement around the levels
- A baseline benchmark
 - No environmental clutters/props
 - The complexity of the level design can be increased to make the benchmark more difficult in the future



э

A D > A P > A D > A D >
Divided into two components



<ロト <回ト < 注ト < 注ト

- Divided into two components
 - Setting up the paths



<ロト <回ト < 注ト < 注ト

Divided into two components

- Setting up the paths
- Placing the NPCs on the level

Divided into two components

- Setting up the paths
- Placing the NPCs on the level
- Images captured using AirSim Python API

- Divided into two components
 - Setting up the paths
 - Placing the NPCs on the level
- Images captured using AirSim Python API
 - RGB

- Divided into two components
 - Setting up the paths
 - Placing the NPCs on the level
- Images captured using AirSim Python API
 - RGB
 - Segmentation



- Divided into two components
 - Setting up the paths
 - Placing the NPCs on the level
- Images captured using AirSim Python API
 - RGB
 - Segmentation
 - Depth Perspective



э

- Divided into two components
 - Setting up the paths
 - Placing the NPCs on the level
- Images captured using AirSim Python API
 - RGB
 - Segmentation
 - Depth Perspective
- Segmentation images used to compute bounding boxes around crowd groups



э

- Divided into two components
 - Setting up the paths
 - Placing the NPCs on the level
- Images captured using AirSim Python API
 - RGB
 - Segmentation
 - Depth Perspective
- Segmentation images used to compute bounding boxes around crowd groups
- Individual bounding boxes processed to obtain a single Ground Truth Bounding Box which should encompass the largest crowd



э

- Divided into two components
 - Setting up the paths
 - Placing the NPCs on the level
- Images captured using AirSim Python API
 - RGB
 - Segmentation
 - Depth Perspective
- Segmentation images used to compute bounding boxes around crowd groups
- Individual bounding boxes processed to obtain a single Ground Truth Bounding Box which should encompass the largest crowd
- The frames are combined into a single video



э

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

- Divided into two components
 - Setting up the paths
 - Placing the NPCs on the level
- Images captured using AirSim Python API
 - RGB
 - Segmentation
 - Depth Perspective
- Segmentation images used to compute bounding boxes around crowd groups
- Individual bounding boxes processed to obtain a single Ground Truth Bounding Box which should encompass the largest crowd
- The frames are combined into a single video
- We will demonstrate the entire process at the end, stay tuned!



э

Time Dilation is employed while capturing images



Time Dilation is employed while capturing images
 Time dilation implies "slowing down" the simulation world



Time Dilation is employed while capturing images

Time dilation implies "slowing down" the simulation world
 1 second elapsed in simulation world clock = n seconds elapsed in real world clock; n > 1



Time Dilation is employed while capturing images

- Time dilation implies "slowing down" the simulation world
- 1 second elapsed in simulation world clock = n seconds elapsed in real world clock; n > 1
- The AirSim API cannot maintain stable frame rates if this is not done



э

Time Dilation is employed while capturing images

- Time dilation implies "slowing down" the simulation world
- 1 second elapsed in simulation world clock = n seconds elapsed in real world clock; n > 1
- The AirSim API cannot maintain stable frame rates if this is not done
- UAV controlled via RC or the API



Time Dilation is employed while capturing images

- Time dilation implies "slowing down" the simulation world
- 1 second elapsed in simulation world clock = n seconds elapsed in real world clock; n > 1
- The AirSim API cannot maintain stable frame rates if this is not done
- UAV controlled via RC or the API
 - Move around the world to capture images of crowd groups



Time Dilation is employed while capturing images

- Time dilation implies "slowing down" the simulation world
- 1 second elapsed in simulation world clock = n seconds elapsed in real world clock; n > 1
- The AirSim API cannot maintain stable frame rates if this is not done
- UAV controlled via RC or the API
 - Move around the world to capture images of crowd groups
 - Possible to maintain stable viewpoint at static height as well



э

Ground Truth Generation Pseudo code

define compute_ground_truth(crowd_bounding_boxes, overlap_threshold, max_iterations):

```
filtered boxes = {}
while current iteration < max iterations:
     foreach box in crowd bounding boxes:
          foreach other box \neq box in crowd bounding boxes:
               if overlap ≥ overlap_threshold:
                    combine box and other box into one box
               else:
                    pick the box with the highest area
                    discard the other box with smaller area
                    add the picked box to filtered boxes
     crowd_bounding_boxes = filtered boxes
     current_iteration += 1
if filtered boxes.length() is 1:
     return filtered boxes[0]
else:
     return the box in filtered boxes with the largest area
```

end compute_ground_truth



・ロト ・ 同ト ・ ヨト ・ ヨト

We hire individuals to annotate the dataset



We hire individuals to annotate the dataset

They should place bounding box around the crowd group they believe to be the largest



ヘロト ヘロト ヘヨト ヘヨト

We hire individuals to annotate the dataset

- They should place bounding box around the crowd group they believe to be the largest
- This is subjective by nature



э

ヘロト ヘロト ヘビト ヘビト

We hire individuals to annotate the dataset

- They should place bounding box around the crowd group they believe to be the largest
- This is subjective by nature
- We measure the Intersection-over-Union of:



э

We hire individuals to annotate the dataset

- They should place bounding box around the crowd group they believe to be the largest
- This is subjective by nature
- ▶ We measure the Intersection-over-Union of:
 - The generated ground truth bounding box of each frame



ヘロト ヘロト ヘビト ヘビト

We hire individuals to annotate the dataset

- They should place bounding box around the crowd group they believe to be the largest
- This is subjective by nature
- ▶ We measure the Intersection-over-Union of:
 - The generated ground truth bounding box of each frame
 - The human annotated bounding box of each frame

э

We hire individuals to annotate the dataset

- They should place bounding box around the crowd group they believe to be the largest
- This is subjective by nature
- We measure the Intersection-over-Union of:
 - The generated ground truth bounding box of each frame
 - The human annotated bounding box of each frame
- Compute 1 IOU (Inverse of IOU / Overlap Error)



We hire individuals to annotate the dataset

- They should place bounding box around the crowd group they believe to be the largest
- This is subjective by nature
- We measure the Intersection-over-Union of:
 - The generated ground truth bounding box of each frame
 - The human annotated bounding box of each frame
- Compute 1 IOU (Inverse of IOU / Overlap Error)
- Compute the average of errors in each sequence of frames



э

Ground Truth Results





<ロト <回ト < 注ト < 注ト

Figure: Good cases where human annotation and computed ground truth agree



э

Ground Truth Results (continued)



Figure: The algorithm picks a larger box, where human annotator believes the person at the bottom should be skipped for a proper box



Ground Truth Results (continued)



Figure: The algorithm fails to pick the proper groups of crowds and ends up selecting a smaller crowd group



イロト イポト イヨト イヨト





A valid question

Vision data based ground truth



A valid question

- Vision data based ground truth
- Erratic at times due to changes in vision properties



э

- A valid question
- Vision data based ground truth
- Erratic at times due to changes in vision properties
- By contrast, a tracking algorithm



э

A valid question

- Vision data based ground truth
- Erratic at times due to changes in vision properties
- By contrast, a tracking algorithm
 - can learn the changes in the patterns



A valid question

- Vision data based ground truth
- Erratic at times due to changes in vision properties
- By contrast, a tracking algorithm
 - can learn the changes in the patterns
 - can also maintain a stable bounding box on the target for a longer time



э
If Ground Truth can be computed, why a tracking algorithm is required?

A valid question

- Vision data based ground truth
- Erratic at times due to changes in vision properties
- By contrast, a tracking algorithm
 - can learn the changes in the patterns
 - can also maintain a stable bounding box on the target for a longer time
 - visual servoing techniques are benefitted from stability of the bounding box



э

(日)

Our Tracking Algorithm: Hybrid Motion Pooling



æ

・ロト ・ 御 ト ・ ヨ ト ・ ヨ ト

 Has two fundamental components: model initializer and optimizer.



э

ヘロト 人間ト 人間ト 人間ト

- Has two fundamental components: model initializer and optimizer.
- Trained model weight: $f = \rho(\chi_{sample})$

- Has two fundamental components: model initializer and optimizer.
- Trained model weight: $f = \rho(\chi_{sample})$
- Training sample pairs: $\chi_{sample} = \{(v_j, c_j)\}_{j=1}^m$



・ロト ・ 一下 ・ ト ・ 日 ・

- Has two fundamental components: model initializer and optimizer.
- Trained model weight: $f = \rho(\chi_{sample})$
- Training sample pairs: $\chi_{sample} = \{(v_j, c_j)\}_{j=1}^m$
- ▶ v_j : extracted visual features. $c_j \in \mathbb{R}^2$: centroid of target.

э

・ロト ・ 一下 ・ ト ・ 日 ・

- Has two fundamental components: model initializer and optimizer.
- Trained model weight: $f = \rho(\chi_{sample})$
- Training sample pairs: $\chi_{sample} = \{(v_j, c_j)\}_{j=1}^m$
- ▶ v_j : extracted visual features. $c_j \in \mathbb{R}^2$: centroid of target.
- Loss is calculated as:

$$\ell(f) = \frac{1}{|\chi_{sample}|} \sum_{j=1}^{m} \|\nabla(v_j * f, c_j)\|^2 + \|\tau f\|^2$$



э

- Has two fundamental components: model initializer and optimizer.
- Trained model weight: $f = \rho(\chi_{sample})$
- Training sample pairs: $\chi_{sample} = \{(v_j, c_j)\}_{j=1}^m$
- ▶ v_j : extracted visual features. $c_j \in \mathbb{R}^2$: centroid of target.
- Loss is calculated as:

$$\ell(f) = rac{1}{|\chi_{sample}|} \sum_{j=1}^{m} \|\nabla(v_j * f, c_j)\|^2 + \|\tau f\|^2$$

f is iteratively updated offline based on m samples.



・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

- Has two fundamental components: model initializer and optimizer.
- Trained model weight: $f = \rho(\chi_{sample})$
- Training sample pairs: $\chi_{sample} = \{(v_j, c_j)\}_{j=1}^m$
- ▶ v_j : extracted visual features. $c_j \in \mathbb{R}^2$: centroid of target.
- Loss is calculated as:

$$\ell\left(f
ight)=rac{1}{\left|\chi_{\textit{sample}}
ight|}\sum_{j=1}^{m}\left\|
abla\left(\textit{v}_{j}*f,\textit{c}_{j}
ight)
ight\|^{2}+\left\| au f
ight\|^{2}$$

- f is iteratively updated offline based on m samples.
- Issue: Can not scale to groups or multiple human interactions.



э

・ロト ・ 『 ト ・ ヨ ト ・ ヨ ト

- Has two fundamental components: model initializer and optimizer.
- Trained model weight: $f = \rho(\chi_{sample})$
- Training sample pairs: $\chi_{sample} = \{(v_j, c_j)\}_{j=1}^m$
- ▶ v_j : extracted visual features. $c_j \in \mathbb{R}^2$: centroid of target.
- Loss is calculated as:

$$\ell\left(f
ight)=rac{1}{\left|\chi_{\textit{sample}}
ight|}\sum_{j=1}^{m}\left\|
abla\left(\textit{v}_{j}*f,\textit{c}_{j}
ight)
ight\|^{2}+\left\| au f
ight\|^{2}$$

- f is iteratively updated offline based on m samples.
- Issue: Can not scale to groups or multiple human interactions.
- Solution: Incorporate Spatio-Temporal Graphs!



 Goal: Capture spatial interaction between humans by building a graph Φ^S.



- Goal: Capture spatial interaction between humans by building a graph Φ^S.
- At time step t: $\gamma_t = \{(\vartheta_t^1, b_t^1), (\vartheta_t^2, b_t^2), ..., (\vartheta_t^n, b_t^{n_t})\}$, where $\vartheta_t^i \in \mathbb{R}^{1 \times d}$



◆日 > < 同 > < 国 > < 国 >

- Goal: Capture spatial interaction between humans by building a graph Φ^S.
- At time step t: $\gamma_t = \{(\vartheta_t^1, b_t^1), (\vartheta_t^2, b_t^2), ..., (\vartheta_t^n, b_t^{n_t})\}$, where $\vartheta_t^i \in \mathbb{R}^{1 \times d}$
- Build graph based on Intersection over Union (IoU) weights:

$$\Phi_{t_{ij}}^{S} = \frac{e^{\zeta_t^{ij}}}{\sum_{j=1}^{n_t} e^{\zeta_t^{ij}}}$$



э

(日)

- Goal: Capture spatial interaction between humans by building a graph Φ^S.
- At time step t: $\gamma_t = \{(\vartheta_t^1, b_t^1), (\vartheta_t^2, b_t^2), ..., (\vartheta_t^n, b_t^{n_t})\}$, where $\vartheta_t^i \in \mathbb{R}^{1 \times d}$
- Build graph based on Intersection over Union (IoU) weights:

$$\Phi^{\mathcal{S}}_{t_{ij}} = rac{e^{\zeta^{ij}_t}}{\sum_{j=1}^{n_t} e^{\zeta^{ij}_t}}$$

• $\Phi_{t_{ij}}^{S} \in \mathbb{R}^{n \times n}$ is spatial adjacency matrix. $\zeta_t^{ij} = IoU(b_t^i, b_t^j)$



э

A D > A D > A D > A D > A

Connect the frames with pair-wise cosine similarity as:

$$\Phi_{t_{ij}}^{\mathcal{T}} = \frac{e^{\cos(\vartheta_t^i, \vartheta_{t+1}^j)}}{\sum_{j=1}^{n_t+1} e^{\cos(\vartheta_t^i, \vartheta_{t+1}^j)}}$$



æ

ヘロト 人間ト 人間ト 人間ト

Connect the frames with pair-wise cosine similarity as:

$$\Phi_{t_{ij}}^{\mathcal{T}} = \frac{e^{\cos(\vartheta_t^i, \vartheta_{t+1}^j)}}{\sum_{j=1}^{n_t+1} e^{\cos(\vartheta_t^i, \vartheta_{t+1}^j)}}$$

• $\Phi_{t_{ij}}^T \in \mathbb{R}^{n \times n+1}$ is temporal adjacency matrix element at (i, j)



ヘロト ヘロト ヘヨト ヘヨト

Connect the frames with pair-wise cosine similarity as:

$$\Phi_{t_{ij}}^{\mathcal{T}} = \frac{e^{\cos(\vartheta_t^i, \vartheta_{t+1}^j)}}{\sum_{j=1}^{n_t+1} e^{\cos(\vartheta_t^i, \vartheta_{t+1}^j)}}$$

Φ^T_{tij} ∈ ℝ^{n×n+1} is temporal adjacency matrix element at (i,j)
 Merge spatial and temporal graph as:

$$\Phi^{ST} = \begin{bmatrix} \Phi_1^S & \Phi_1^T & 0 & \cdots & 0\\ 0 & \Phi_2^S & \Phi_2^T & \cdots & 0\\ 0 & 0 & \Phi_3^S & \cdots & 0\\ \vdots & \vdots & \vdots & \ddots & 0\\ 0 & 0 & 0 & \cdots & \Phi_t^S \end{bmatrix}$$



Connect the frames with pair-wise cosine similarity as:

$$\Phi_{t_{ij}}^{T} = \frac{e^{\cos(\vartheta_t^i, \vartheta_{t+1}^j)}}{\sum_{j=1}^{n_t+1} e^{\cos(\vartheta_t^i, \vartheta_{t+1}^j)}}$$

Φ^T_{tij} ∈ ℝ^{n×n+1} is temporal adjacency matrix element at (i,j)
 Merge spatial and temporal graph as:

$$\Phi^{ST} = \begin{bmatrix} \Phi_1^S & \Phi_1^T & 0 & \cdots & 0 \\ 0 & \Phi_2^S & \Phi_2^T & \cdots & 0 \\ 0 & 0 & \Phi_3^S & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & 0 \\ 0 & 0 & 0 & \cdots & \Phi_t^S \end{bmatrix}$$

• $\Phi^{ST} \in \mathbb{R}^{N \times N}$ is composed of Φ_t^S and Φ_t^T



Connect the frames with pair-wise cosine similarity as:

$$\Phi_{t_{ij}}^{T} = \frac{e^{\cos(\vartheta_t^i, \vartheta_{t+1}^j)}}{\sum_{j=1}^{n_t+1} e^{\cos(\vartheta_t^i, \vartheta_{t+1}^j)}}$$

Φ^T_{tij} ∈ ℝ^{n×n+1} is temporal adjacency matrix element at (i, j)
 Merge spatial and temporal graph as:

$$\Phi^{ST} = \begin{bmatrix} \Phi_1^S & \Phi_1^T & 0 & \cdots & 0\\ 0 & \Phi_2^S & \Phi_2^T & \cdots & 0\\ 0 & 0 & \Phi_3^S & \cdots & 0\\ \vdots & \vdots & \vdots & \ddots & 0\\ 0 & 0 & 0 & \cdots & \Phi_t^S \end{bmatrix}$$

• $\Phi^{ST} \in \mathbb{R}^{N \times N}$ is composed of Φ_t^S and Φ_t^T • Update using GCN and get $\Lambda \in \mathbb{R}^{T \times dm}$:

$$\mathbb{C}^{\prime+1} = \sigma(\mathbb{C}^{\prime} + \lambda^{-\frac{1}{2}} \Phi^{ST} \lambda^{-\frac{1}{2}} \mathbb{C}^{\prime} W^{\prime})$$



Combine Λ with f from DiMP at each timestep using low-rank bilinear pooling as:

$$\beta_i = \sigma \left(W_f f \right) \odot \sigma \left(W_{\Lambda} \Lambda_i \right)$$

æ

ヘロト ヘロト ヘヨト ヘヨト

Combine Λ with f from DiMP at each timestep using low-rank bilinear pooling as:

$$\beta_i = \sigma \left(W_f f \right) \odot \sigma \left(W_{\Lambda} \Lambda_i \right)$$

• Calculate final filter ∂_T by performing temporal max-pooling:

$$\partial_{T} = MaxPool_{T} ([\beta_{1}, \beta_{2}, \dots, \beta_{T-1}]) \in \mathbb{R}^{1 \times D_{\Lambda}}$$



э

Combine Λ with f from DiMP at each timestep using low-rank bilinear pooling as:

$$\beta_i = \sigma(W_f f) \odot \sigma(W_{\Lambda} \Lambda_i)$$

Calculate final filter \u03c8_T by performing temporal max-pooling:

$$\partial_{T} = MaxPool_{T} ([\beta_{1}, \beta_{2}, \dots, \beta_{T-1}]) \in \mathbb{R}^{1 \times D_{\Lambda}}$$

• dimension of D_{Λ} does not depend on training video length T



э

Combine Λ with f from DiMP at each timestep using low-rank bilinear pooling as:

$$\beta_i = \sigma \left(W_f f \right) \odot \sigma \left(W_{\Lambda} \Lambda_i \right)$$

Calculate final filter \u03c8_T by performing temporal max-pooling:

$$\partial_{T} = MaxPool_{T} ([\beta_{1}, \beta_{2}, \dots, \beta_{T-1}]) \in \mathbb{R}^{1 \times D_{\Lambda}}$$

dimension of D_Λ does not depend on training video length T
 T: temporal information for HyMP. m: for DiMP



э

▶ Initially, pre-train the DiMP weight model f from the UUCT dataset for T = m = 10.

- ▶ Initially, pre-train the DiMP weight model f from the UUCT dataset for T = m = 10.
- For test window pairs, loss is calculated as:

$$\ell_{t \operatorname{arg} et} = \frac{1}{\mathbb{N}} \sum_{i=1}^{\mathbb{N}} \sum_{j=1}^{m} \left\| \kappa \left(\mathsf{v}_{j} \ast \partial_{T}^{i}, \mathsf{g}_{\mathsf{c}_{j}} \right) \right\|^{2}$$



э

ヘロト ヘ週ト ヘヨト ヘヨト

- Initially, pre-train the DiMP weight model f from the UUCT dataset for T = m = 10.
- For test window pairs, loss is calculated as:

$$\ell_{t \text{ arg } et} = rac{1}{\mathbb{N}} \sum_{i=1}^{\mathbb{N}} \sum_{j=1}^{m} \left\| \kappa \left(\mathsf{v}_{j} * \partial_{T}^{i}, \mathsf{g}_{\mathsf{c}_{j}}
ight)
ight\|^{2}$$

$$\kappa(v,c) = \left\{ egin{array}{ll} v-c, & ext{if } c > arepsilon \ \max(0,c), & ext{otherwise} \end{array}
ight.$$

 $\blacktriangleright \ \ell_{final} = \eta \ell_{t \, arg \, et} + \ell_{bb}$



э

(日) (圖) (目) (日) (日)

Initially, pre-train the DiMP weight model f from the UUCT dataset for T = m = 10.

For test window pairs, loss is calculated as:

$$\ell_{t \operatorname{arg} et} = \frac{1}{\mathbb{N}} \sum_{i=1}^{\mathbb{N}} \sum_{j=1}^{m} \left\| \kappa \left(\mathsf{v}_{j} \ast \partial_{T}^{i}, \mathsf{g}_{c_{j}} \right) \right\|^{2}$$

$$\kappa\left(\mathbf{v},\mathbf{c}
ight) = \left\{egin{array}{ll} \mathbf{v}-\mathbf{c}, & ext{if } \mathbf{c} > arepsilon \ \max(\mathbf{0},\mathbf{c}), & ext{otherwise} \end{array}
ight.$$



э

・ロト ・四ト ・ヨト ・ヨト

Online Tracking and Algorithm

	Algorithm 1 Online HyMP Tracking Algorithm							
	Input : Initial Trained Model weight ∂_T , Initial Test sample $\chi_{test} =$							
	(v_0, c_0) . New frames v_q where $q > 0$							
	Output: Estimated new positions with updated model							
1	$\chi_{sample} \leftarrow augment \chi_{test} = (v_0, c_0)$ pair into 10 samples							
2	$\gamma_0 \leftarrow$ apply detector and fetch feature, box pair from χ_{test}							
3	$\gamma \leftarrow repeat \ \gamma_0 \ 10 \ times$							
4	$\partial_T \leftarrow$ forward pass χ_{sample} and γ to the model ρ and perform 10 training iterations							
5	repeat							
	/* bounding box prediction stage */							
6	$y^{tc} \in \mathbb{R}^2 \leftarrow \text{Perform convolution on } \delta(v_q) * \partial_T$							
7	$y^{bb} \in \mathbb{R}^4, p_a \leftarrow \text{Regress candidate bounding boxes and choose}$							
	one with the highest confidence score.							
8	plot y^{bb} in the frame v_q							
	/* model update step */							
9	if $p_q > \varepsilon$ then							
	<pre>/* append frame (+) to memory */</pre>							
10	$\chi_{sample} \leftarrow \chi_{sample} + \delta(v_q)$							
11	$\gamma_t \leftarrow \text{feature, box pair for } v_q$							
12	$\gamma \leftarrow \gamma - \gamma_t$							
13	end							
14	if $length(\chi_{sample}) > 50$ then							
	/* remove frame (-) from memory */							
15	$\chi_{sample} \leftarrow \chi_{sample} - v_0$							
16	$\gamma \leftarrow \gamma + \gamma_0$							
17	end							
18	it 30 new trames are appended then							
19	$\sigma_T \leftarrow$ forward pass χ_{sample} and γ to the model ρ and perform							
20	end							
21 until Until novel v. without bounding boxes is received								
41	until onthi nover v _q without bounding boxes is received							



ヘロト 人間ト 人間ト 人間ト

Experimental Evaluation

Table: Comparison of algorithms' AUC on the benchmark datasets. For VOT2018, we compared EAO score [KLM⁺18]. HyMP-X is HyMP with Xception [Cho17] as backbone feature extractor

Tracker	VOT-2018	OTB100	UAV123	UAV20L	DTB70	UUCT		
Correlation Filter Trackers								
Staple [BVG ⁺ 16]	16.9	58.6	45.0	33.1	35.1	19.6		
SRDCF [DHSKF15]	11.9	59.8	46.4	34.3	36.3	20.3		
STRCF [LTZ ⁺ 18]	14.3	64.1	48.1	35.4	40.7	20.3		
Tracekrs applied Correlation Filter on deeply learned features								
HCFT [MHYY18b]	19.9	64.7	48.6	36.8	41.5	22.7		
LCT [MHYY18a]	22.1	65.2	49.4	35.9	43.1	25.5		
Tracekrs exploited end-to-end deep learning pipeline								
GCT [GZX19]	27.6	64.8	50.8	46.1	44.2	27.2		
UPDT [BJD ⁺ 18]	38.3	70.2	54.5	49.5	45.7	28.1		
DiMP-18 [BDGT19]	40.2	66	64.3	51.7	46.9	29.8		
DiMP-50 [BDGT19]	43.1	68.4	65.4	52.8	47.5	30.6		
End-to-end HyMP Tracker (ours)								
HyMP-18	40.1	66.8	66.2	52.5	47.1	31.7		
HyMP-50	42.8	69.1	67.6	52.9	48.6	32.9		
HyMP-X	43.6	69.4	68.2	53.7	49.2	34.0		



ヘロア 人間 アメヨア メヨア

୍ର ୧୦୦ ୯ ୦

э

Algorithm Performance Visualization on UUCT Dataset



Figure: Overall benchmark of SOTA trackers on UUCT



Figure: Benchmark on Crowd Split and Dispersed Out of View (CSDOV)



Thank You!



◆□ ▶ ◆圖 ▶ ◆ 壹 ▶ ◆ 壹 ▶

References



Goutam Bhat, Martin Danelljan, Luc Van Gool, and Radu Timofte, *Learning discriminative model prediction for tracking*, Proceedings of the IEEE International Conference on Computer Vision, 2019.



Goutam Bhat, Joakim Johnander, Martin Danelljan, Fahad Shahbaz Khan, and Michael Felsberg, *Unveiling the power of deep tracking*, Proceedings of the European Conference on Computer Vision (ECCV), 2018.



Luca Bertinetto, Jack Valmadre, Stuart Golodetz, Ondrej Miksik, and Philip HS Torr, *Staple: Complementary learners for real-time tracking*, Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 1401–1409.



François Chollet, *Xception: Deep learning with depthwise separable convolutions*, Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 1251–1258.



Martin Danelljan, Gustav Hager, Fahad Shahbaz Khan, and Michael Felsberg, *Learning spatially regularized correlation filters for visual tracking*, Proceedings of the IEEE international conference on computer vision, 2015.



Junyu Gao, Tianzhu Zhang, and Changsheng Xu, *Graph convolutional tracking*, Proceedings of the IEEE conference on computer vision and pattern recognition, 2019, pp. 4649–4659.



Matej Kristan, Ales Leonardis, Jiri Matas, Michael Felsberg, Roman Pflugfelder, Luka ^{*}Cehovin Zajc, Tomas Vojir, Goutam Bhat, Alan Lukezic, Abdelrahman Eldesokey, et al., *The sixth visual object tracking vot2018 challenge results*, Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 0–0.



Feng Li, Cheng Tian, Wangmeng Zuo, Lei Zhang, and Ming-Hsuan Yang, *Learning spatial-temporal regularized correlation filters for visual tracking*, Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018.



1

Chao Ma, Jia-Bin Huang, Xiaokang Yang, and Ming-Hsuan Yang, Adaptive correlation filters with long-term and short-term memory for object tracking, International Journal of Computer Vision 126 (2018), no. 8.



_____, Robust visual tracking via hierarchical convolutional features, IEEE transactions on pattern analysis and machine intelligence (2018) (en).