# Supplementary Material of "Global2Local: Efficient Structure Search for Video Action Segmentation"

Shang-Hua Gao[1*]     Qi Han[1*]     Zhong-Yu Li[1]
Pai Peng[2]     Liang Wang[3]     Ming-Ming Cheng[1] [†]
TKLNDST, CS, Nankai University[1]     Tencent[2]     NLPR[3]

http://mmcheng.net/g2lsearch

| GTEA | F@0.1 | F@0.25 | F@0.5 | Edit | Acc |
|------|-------|--------|-------|------|-----|
| *Spatial CNN* | 41.8 | 36.0 | 25.1 | - | 54.1 |
| *Bi-LSTM* | 66.5 | 59.0 | 43.6 | - | 55.5 |
| *Dilated TCN* | 58.8 | 52.2 | 42.2 | - | 58.3 |
| *ST-CNN* | 58.7 | 54.4 | 41.9 | - | 60.6 |
| *TUnet* | 67.1 | 63.7 | 51.9 | 60.3 | 59.9 |
| *ED-TCN* | 72.2 | 69.3 | 56.0 | - | 64.0 |
| *TResNet* | 74.1 | 69.9 | 57.6 | 64.4 | 65.8 |
| *TricorNet* | 76.0 | 71.1 | 59.2 | - | 64.8 |
| *TRN* | 77.4 | 71.3 | 59.1 | 72.2 | 67.8 |
| *TDRN* | 79.2 | 74.4 | 62.7 | 74.1 | 70.1 |
| *MS-TCN* | 87.5 | 85.4 | 74.6 | 81.4 | 79.2 |
| *Reproduce* | 87.1 | 83.6 | 70.4 | 81.1 | 75.5 |
| *Ours-MS-TCN* | 89.9 | 87.3 | 75.8 | 84.6 | 78.5 |

Table 1. Comparison with state-of-the-art methods on the GTEA dataset.

## 1. Performance on the GTEA dataset

We also compare our proposed global-to-local search with existing action segmentation methods on the small scale GTEA dataset, as shown in Tab. 1. Based on the MS-TCN architecture, the global-to-local searched structure surpasses the human-designed baseline with 2.4% on F@0.1. Also, our global-to-local based MS-TCN has a considerable performance gain compared with existing methods.

## 2. Searching Cost

We report the cost of our proposed global-to-local search method. When cooperating with MS-TCN, the size of the receptive field combination search space is $1024^{40}$. The cost of searching on such a huge space using existing methods is unaffordable. Our proposed global-to-local search decomposes the searching process into the global and local

---

*Equal contribution
[†]M.M. Cheng (cmm@nankai.edu.cn) is the corresponding author.

| GPU Hours | BreakFast | 50Salads | GTEA |
|-----------|-----------|----------|------|
| *Global Search* | 144h | 9h | 1h |
| *Local Search* | 2.2h | 0.15h | 0.05h |

Table 2. GPU hours of the global and local search on each fold of different datasets using the RTX 2080Ti GPU.

| Actual Searching Hours | BreakFast | 50Salads | GTEA |
|------------------------|-----------|----------|------|
| *Global Search* (8 GPU) | 42h | 11h | 5h |
| *Local Search* (1 GPU) | 3.3h | 0.2h | 0.07h |

Table 3. The actual searching time of the global and local search on each fold of different datasets using the RTX 2080Ti GPU×8 server.

search to find the combination in a coarse-to-fine manner.

In the global search process, we search structures on each fold of datasets. Since the main bottleneck of the search method is the GPU resources, we report the GPU hours of the proposed global-to-local search in Tab. 2. In our work, 1 GPU hour means using one RTX 2080Ti GPU for an hour. For each fold in the dataset, the global search takes 144 GPU hours on the BreakFast dataset, 9 GPU hours on the 50Salads dataset, and 1 GPU hours on the GTEA dataset. Through the global search, we can find multiple new well-performed structures that have different patterns and achieve better performance than human-designed patterns. The local search further fine-tunes the global searched structures in the dense but local search space. Based on one of the global-searched structures, the local search takes about 2.2 GPU hours on the BreakFast dataset, 0.15 GPU hours on the 50Salads dataset, and 0.05 hours on the GTEA dataset.

Due to the CPU and disk IO speed limitation in actual experiments, the search time is longer than the GPU hour. We report the actual searching time on an RTX 2080Ti GPU×8 server in Tab. 3. For each fold of the dataset, the global search takes 42 hours on the BreakFast dataset, 11 hours on

the 50Salads dataset, and 5 hours on the GTEA dataset on an RTX 2080Ti GPU×8 server. The local search takes 3.3 hours on the BreakFast dataset, 12 minutes on the 50Salads dataset, and 4 minutes on the GTEA dataset, with a single RTX 2080Ti GPU.

## 3. Common Patterns in the Searched Structures

The global search objective is to find more well-performed receptive field combinations that have different patterns than human-designings. We visualize the top-5 well-performed searched structures on each fold of different datasets. The visualization of searched structures of the BreakFast dataset is shown in Fig. 1, Fig. 2, Fig. 3, and Fig. 4. The visualization of searched structures of the 50salads dataset is shown in Fig. 5, Fig. 6, Fig. 7, Fig. 8, and Fig. 9. The visualization of searched structures of the GTEA dataset is shown in Fig. 10, Fig. 11, Fig. 12, and Fig. 13. We visualize the common receptive field combination patterns in the well-performed searched structures in Tab. 4. Common receptive field combination patterns vary among different datasets. Also, different stages in the MS-TCN have different common patterns.

Table 4. Common receptive field combination patterns in the well-performed searched structures.

| Datasets | stage1 | stage2 | stage3 | stage4 |
|---|---|---|---|---|
| *BreakFast* | 512 64 256 | 4 1 64 2 1024 8 256 | 4 32 8 2 2 8 8 | 1024 16 1 16 128 |
| | 2 128 512 32 1024 8 64 | 64 512 1024 8 1024 128 | 2 32 8 4 | 128 1024 128 64 128 |
| | 2 128 512 64 256 32 | 8 4 1 16 512 1024 | 32 8 2 2 | 1 16 128 1024 128 64 128 |
| *50Salads* | 32 512 1024 | 4 1 64 2 1024 | 1 4 64 32 4 8 1 | 512 16 128 2 32 1 4 |
| | 32 1024 128 32 512 1024 | 1 64 2 1024 8 256 | 256 1 1 4 64 | 128 2 32 1 |
| | 2 512 512 1024 256 1024 256 | 4 1 64 2 1024 8 256 | 128 4 64 32 | 2 32 1 4 64 128 |
| *GTEA* | 512 64 256 512 16 | 512 32 2 1 4 2 16 | 32 4 128 1 16 16 512 | 4 8 2 32 16 |
| | 1 2 1 32 512 64 256 | 1 1 32 32 16 | 1 256 16 512 4 | 32 8 16 1 1024 32 |
| | 512 8 256 64 16 | 1 64 1 1 32 | 8 128 1 16 | 2 32 16 1 512 64 32 |



Figure 1. Well-performed structures on fold 1 of the BreakFast dataset.



Figure 2. Well-performed structures on fold 2 of the BreakFast dataset.



Figure 3. Well-performed structures on fold 3 of the BreakFast dataset.

Figure 4. Well-performed structures on fold 4 of the BreakFast dataset.


Figure 5. Well-performed structures on fold 1 of the 50Salads dataset.
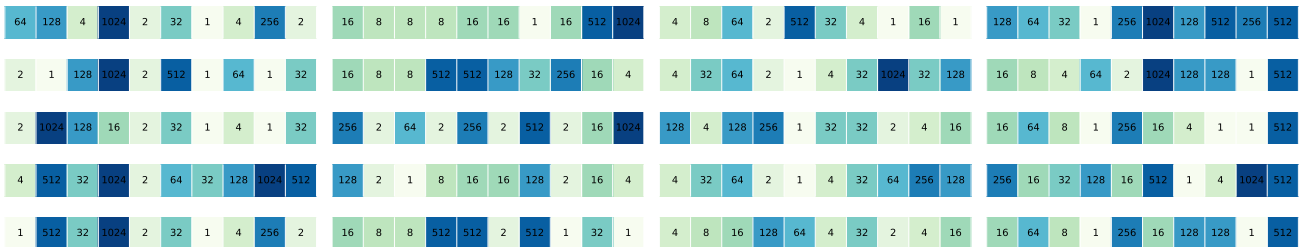

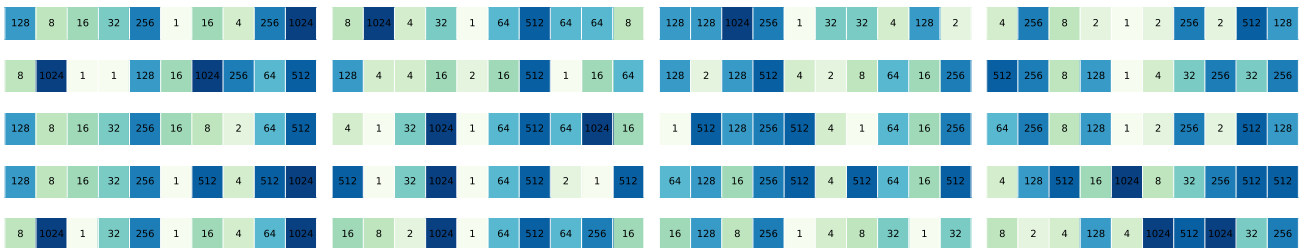Figure 6. Well-performed structures on fold 2 of the 50Salads dataset.


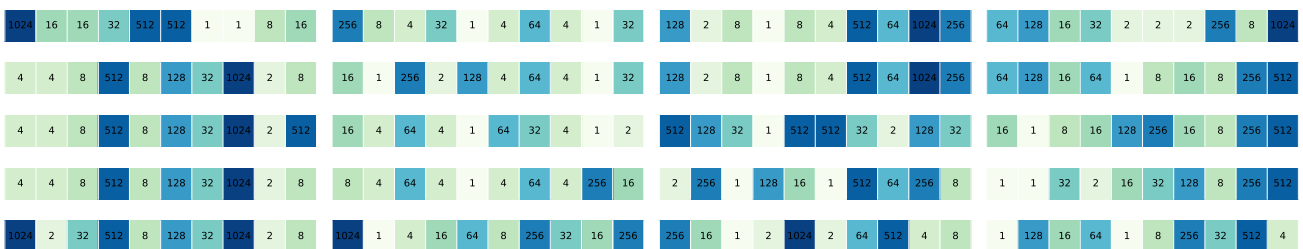Figure 7. Well-performed structures on fold 3 of the 50Salads dataset.


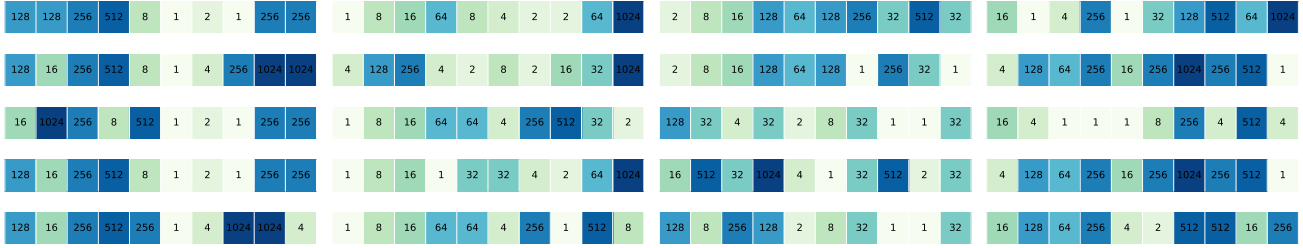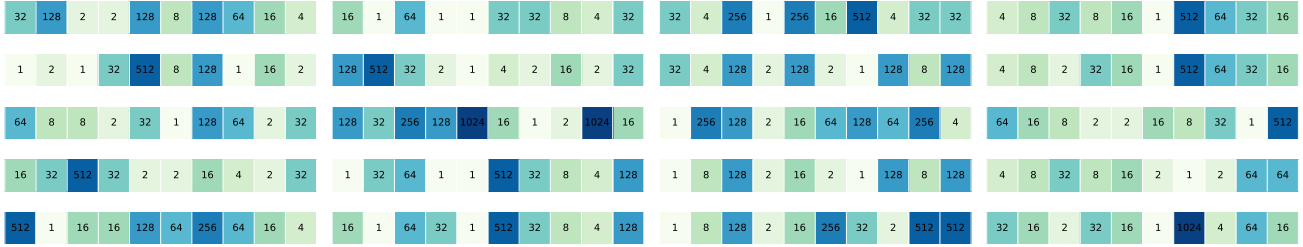Figure 8. Well-performed structures on fold 4 of the 50Salads dataset.

Figure 9. Well-performed structures on fold 5 of the 50Salads dataset.

Figure 10. Well-performed structures on fold 1 of the GTEA dataset.
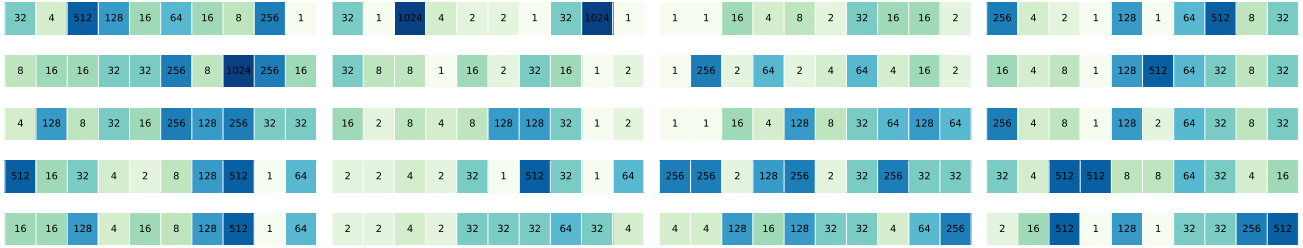
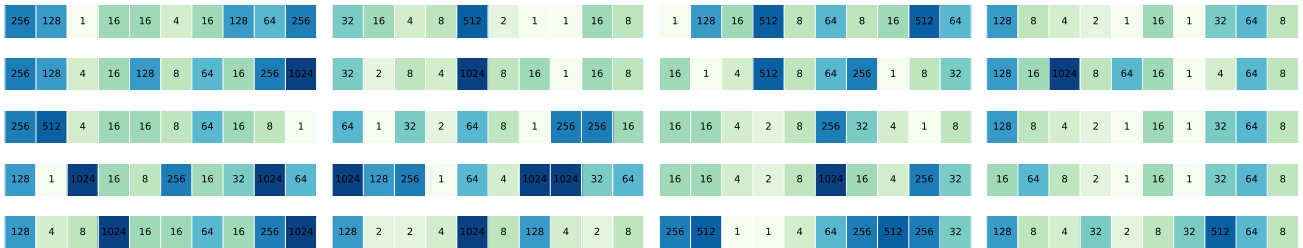Figure 11. Well-performed structures on fold 2 of the GTEA dataset.

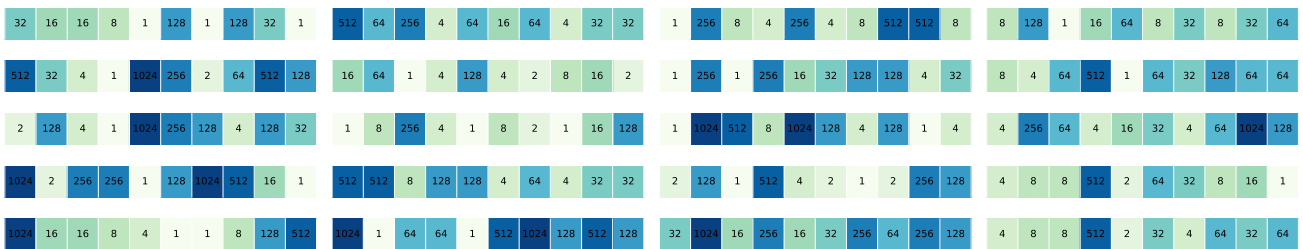Figure 12. Well-performed structures on fold 3 of the GTEA dataset.

Figure 13. Well-performed structures on fold 4 of the GTEA dataset.