

Fast Video Object Segmentation using the Global Context Module

Yu Li^{1*}, Zhuoran Shen^{2*}, and Ying Shan¹

¹ Applied Research Center (ARC), Tencent PCG

² The University of Hong Kong

1 Proof of Equivalence between the Global Context Module and the Space-Time Memory Module

1.1 Mathematical Formulations of the Global Context Module and the Space-Time Memory Module

The global context (GC) module consists of three steps, namely context extraction, context update, and context distribution. The context extraction step can be expressed as

$$\mathbf{C}_t = k(\mathbf{X}_t)^\top v(\mathbf{X}_t), \quad (1)$$

where t is the index of the current frame, \mathbf{X}_t is the input feature, \mathbf{C}_t is the global context feature of the current frame, and k, v are the key and value generation functions. The context update step is

$$\mathbf{G}_t = \frac{t-1}{t}\mathbf{G}_{t-1} + \frac{1}{t}\mathbf{C}_t, \quad (2)$$

where \mathbf{G}_t is the global context feature of the first t frames. The context distribution step is

$$\mathbf{R}_t = q(\mathbf{X}_t)\mathbf{G}_{t-1}, \quad (3)$$

where \mathbf{R}_t is the distributed global context feature and q is the query generation function.

The space-time memory (STM) [1] module similarly consists of three steps, namely memory production, memory write, and memory read. The memory production step can be expressed as

$$\begin{aligned} \mathbf{K}_{P,t} &= k(\mathbf{X}_t), \\ \mathbf{V}_{P,t} &= v(\mathbf{X}_t), \end{aligned} \quad (4)$$

where t is the index of the current frame, \mathbf{X}_t is the input feature, $\mathbf{K}_{M,t}, \mathbf{V}_{M,t}$ are the key and value produced from the current frame, and k, v are the key and value generation functions. The memory write operation is

$$\begin{aligned} \mathbf{K}_{M,t} &= \mathbf{K}_{M,t-1} \odot \mathbf{K}_{P,t}, \\ \mathbf{V}_{M,t} &= \mathbf{V}_{M,t-1} \odot \mathbf{V}_{P,t}, \end{aligned} \quad (5)$$

* Both authors contributed equally. This work was done when Zhuoran was interning at Tencent.

where $\mathbf{K}_{M,t}, \mathbf{V}_{M,t}$ are the key and value in the memory from the first t frames and \odot denotes concatenation along the spatial dimension. The memory read step is

$$\mathbf{E}_t = \frac{1}{t} (q(\mathbf{X}_t) \mathbf{K}_{M,t-1}^\top) \mathbf{V}_{M,t-1}, \quad (6)$$

where \mathbf{E}_t is the memory reading and q is the query generation function.

1.2 Proof of Equivalence

Expanding Equation (2) by substituting it to itself gives

$$\mathbf{G}_{t-1} = \frac{1}{t-1} \sum_{f=1}^{t-1} \mathbf{C}_f. \quad (7)$$

Note that \mathbf{G}_0 is the zero matrix. Substituting Equation (7) into Equation (3) results in

$$\mathbf{R}_t = \frac{1}{t-1} q(\mathbf{X}_t) \sum_{f=1}^{t-1} \mathbf{C}_f. \quad (8)$$

Combining Equations (1) and (8), we have

$$\mathbf{R}_t = \frac{1}{t-1} q(\mathbf{X}_t) \sum_{f=1}^{t-1} k(\mathbf{X}_f)^\top v(\mathbf{X}_f). \quad (9)$$

Similarly, expanding Equations (5) gives

$$\begin{aligned} \mathbf{K}_{M,t-1} &= \bigodot_{f=1}^{t-1} \mathbf{K}_{P,f}, \\ \mathbf{V}_{M,t-1} &= \bigodot_{f=1}^{t-1} \mathbf{V}_{P,f}, \end{aligned} \quad (10)$$

where \bigodot stands for concatenation along the spatial dimension. Substituting Equations (4) into Equation (10) results in

$$\begin{aligned} \mathbf{K}_{M,t-1} &= \bigodot_{f=1}^{t-1} k(\mathbf{X}_f), \\ \mathbf{V}_{M,t-1} &= \bigodot_{f=1}^{t-1} v(\mathbf{X}_f). \end{aligned} \quad (11)$$

Combining Equations (6) and (11), we have

$$\begin{aligned}
\mathbf{E}_t &= \frac{1}{t-1} \left(q(\mathbf{X}_t) \bigcirc_{f=1}^{t-1} k(\mathbf{X}_f)^\top \right) \bigcirc_{f=1}^{t-1} v(\mathbf{X}_f) \\
&= \frac{1}{t-1} \left(\bigcirc_{f=1}^{t-1} q(\mathbf{X}_t) k(\mathbf{X}_f)^\top \right) \bigcirc_{f=1}^{t-1} v(\mathbf{X}_f) \\
&= \frac{1}{t-1} \sum_{f=1}^{t-1} q(\mathbf{X}_t) k(\mathbf{X}_f)^\top v(\mathbf{X}_f) \\
&= \frac{1}{t-1} q(\mathbf{X}_t) \sum_{f=1}^{t-1} k(\mathbf{X}_f)^\top v(\mathbf{X}_f)
\end{aligned} \tag{12}$$

Comparing Equations (9) and (12), we have

$$\mathbf{R}_t = \mathbf{E}_t, \tag{13}$$

which completes the proof.

1.3 Impact of Softmax Normalization

In practice, in addition to the operations in (6), the STM module usually uses a softmax along the rows on $q(\mathbf{X}_t) \mathbf{K}_{M,t-1}^\top$. This operation makes each row sum up to 1. The interpretation of this normalization is that it makes each row, which corresponds to a query pixel, represent a probabilistic distribution of attention over all locations in the memory.

Since softmax is non-linear, the GC module cannot perfectly recreate its effect. However, the GC module applies two softmax operations, one on each row of $q(\mathbf{X}_t)$ and the other on each column of $k(\mathbf{X}_t)$. Then, if one hypothetically multiplies these two normalized matrices, the product matrix will share the critical property of $\text{softmax}(q(\mathbf{X}_t) \mathbf{K}_{M,t-1}^\top)$, i.e. each row sums up to 1 and represents a probabilistic distribution of attention for the query pixel over all locations in the past frames. Therefore, the two operations together are an effective approximation of the single softmax in STM. Thus, the GC module remains approximately equivalent to the STM module even with the presence of the softmax operations.

References

1. Oh, S.W., Lee, J.Y., Xu, N., Kim, S.J.: Video object segmentation using space-time memory networks. In: ICCV (2019)