

# Spatio-Temporal Self-Organizing Map Deep Network for Dynamic Object Detection from Videos

Yang Du, Chunfeng Yuan, Bing Li, Weiming Hu and Stephen Maybank

National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Science, Beijing, China  
 duyang2014@ia.ac.cn, cfyuan@nlpr.ia.ac.cn, bli@nlpr.ia.ac.cn, wmu@nlpr.ia.ac.cn, sjmaybank@dcs.bbk.ac.uk

## Introduction:

- Modern detection algorithms are generally attained by background modeling. The difficulty of background modeling is to tackle background motions.

## Motivations:

- We conclude that the motions of complex background mainly have two properties:
  - Variation of the global background in the space. It is mainly caused by the zoom, translation, jitter, etc, of the camera. We refer to it as the spatial property of background motion.
  - Variation of the local background in the time. It mainly indicates the dynamic elements in background and at different frames, such as river, fountain and bad weather. We refer to it as the temporal property of background motion.

## Contributions:

- We propose a new STSOM and train it in two aspects, using the whole frames from spatial perspective and using the sequence of a pixel over time from temporal perspective. Then we propose a new method based on Bayesian parameter estimation to automatically learn the spatio-temporal threshold of background filtering. In order to further accurately model the complex background, we stack multiple STSOMs to form a deep network with a STSOM as a layer. The different parts of complex background are accurately modeled by different layers. The dynamic objects are detected by filtering out the background layer by layer and the segments are increasingly accurate with the deeper layer.

## Self-Organizing Map:

- $F(t)$ : the t-th input
- $w(t)$ : weight vector
- $c$ : the index of the winner node, which is the one with the weight vector  $w(t)$  that has the smallest Euclidean distance from  $F(t)$

$$c = \operatorname{argmin}_q \{ \|F(t) - w_q(t)\| \}$$

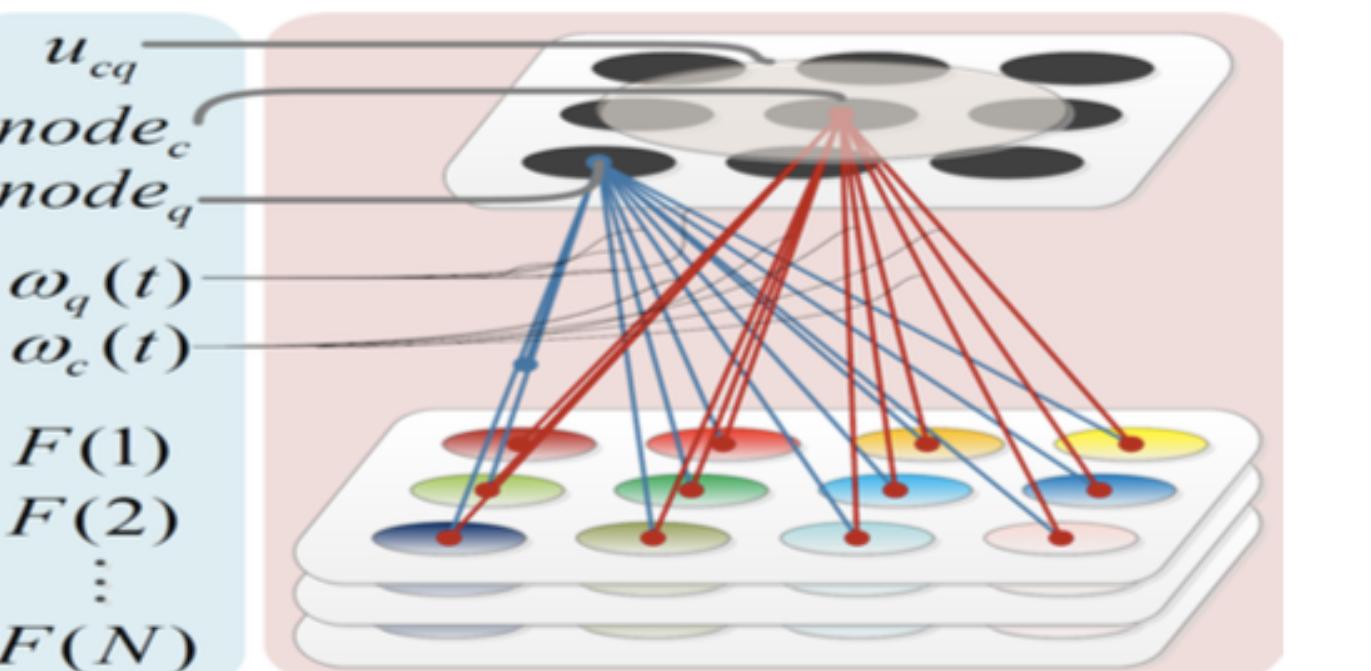


Figure 1. Structure of a self-organizing map.

- neighborhood function:

$$u_{cq} = \exp\left(-\frac{\|q-c\|}{2\sigma^2}\right)$$

- Learning rule:

$$w_q(t+1) = w_q(t) + u_{cq} * \alpha(t) * [F(t) - w_q(t)]$$

## STSOM Deep Network for Dynamic Object Detection :

- A new frame of a video enters into this network and then its dynamic objects will be extracted layer by layer till the last output layer.

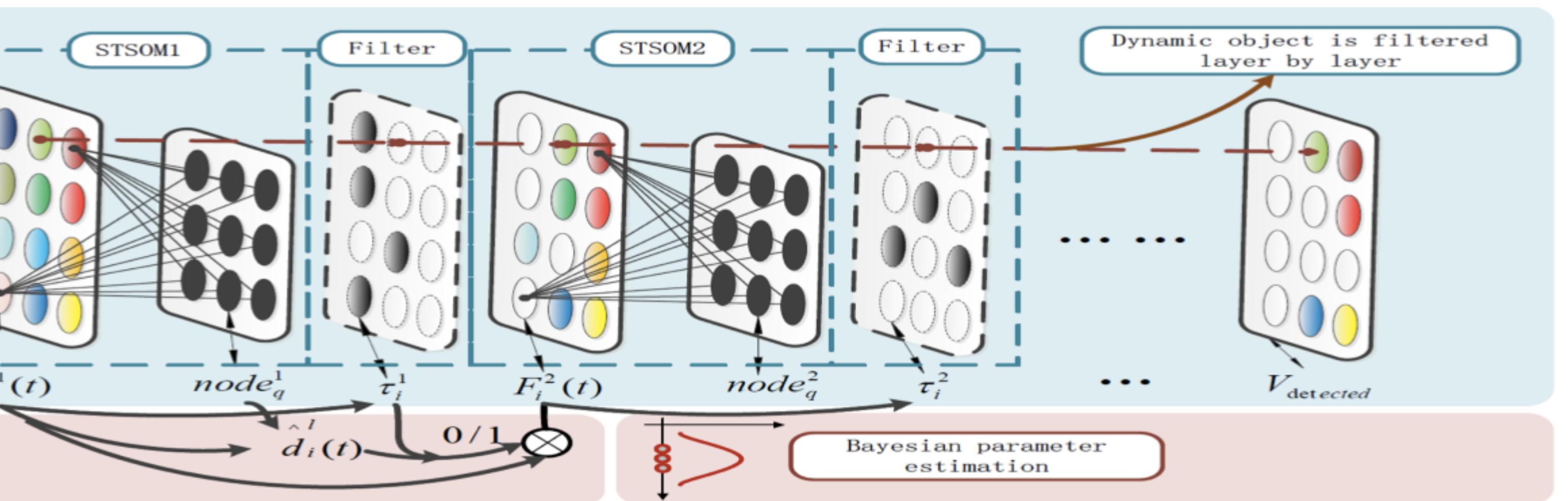


Figure 2. The architecture of the STSOM deep network.

### Spatial Weight Updating

$$D_{spatial,q}^l(t) = \sum_{i=1}^P d_{iq}^l(t)$$

$$q_{spatial}^*(t) = \operatorname{argmin}_q \{ D_{spatial,q}^l(t), q = 1, 2, \dots, Q \}$$

$$w_{iq}^l(t+1) = u_q * (w_{iq}^l(t) + \alpha_{train} * (F_i^l(t) - w_{iq}^l(t)))$$

### Temporal Weight Updating

$$q_{temporal}^*(t) = \operatorname{argmin}_q \{ d_{iq}^l(t), q = 1, 2, \dots, Q \}$$

### Forward Propagation

$$\hat{\mu}_i = \frac{N/\sigma_i^2}{N/\sigma_i^2 + 1/\hat{\delta}_i^2} F_{ave,i}^l + \frac{1/\hat{\delta}_i^2}{N/\sigma_i^2 + 1/\hat{\delta}_i^2} \hat{\gamma}_i$$

$$L(\hat{\mu}_i, \sigma_i^2) = \sum_{t=1}^N \ln N(F_i^l(t) | \hat{\mu}_i, \sigma_i^2).$$

$$\hat{\mu}_i|_{\partial \hat{\mu}_i=0} = \operatorname{argmax}_{\hat{\mu}_i} \{ L(\hat{\mu}_i, \sigma_i^2) \}$$

$$\tau_{spatial,i}^l = \max \{ D_{spatial,q}^l, q = 1, 2, \dots, Q \} / P$$

$$\tau_{temporal,i}^l = \max \{ d_{iq}^l, q = 1, 2, \dots, Q \}$$

$$d_{iq}^l(t) = \|(v_i^l(t)s_i^l(t) \cos(h_i^l(t)), v_i^l(t)s_i^l(t) \sin(h_i^l(t)), \\ v_i^l(t)) - (w_{v_{iq}}^l(t)w_{s_{iq}}^l(t) \cos(w_{h_{iq}}^l(t)), \\ w_{v_{iq}}^l(t)w_{s_{iq}}^l(t) \sin(w_{h_{iq}}^l(t)), w_{v_{iq}}^l(t))\|_2^2.$$

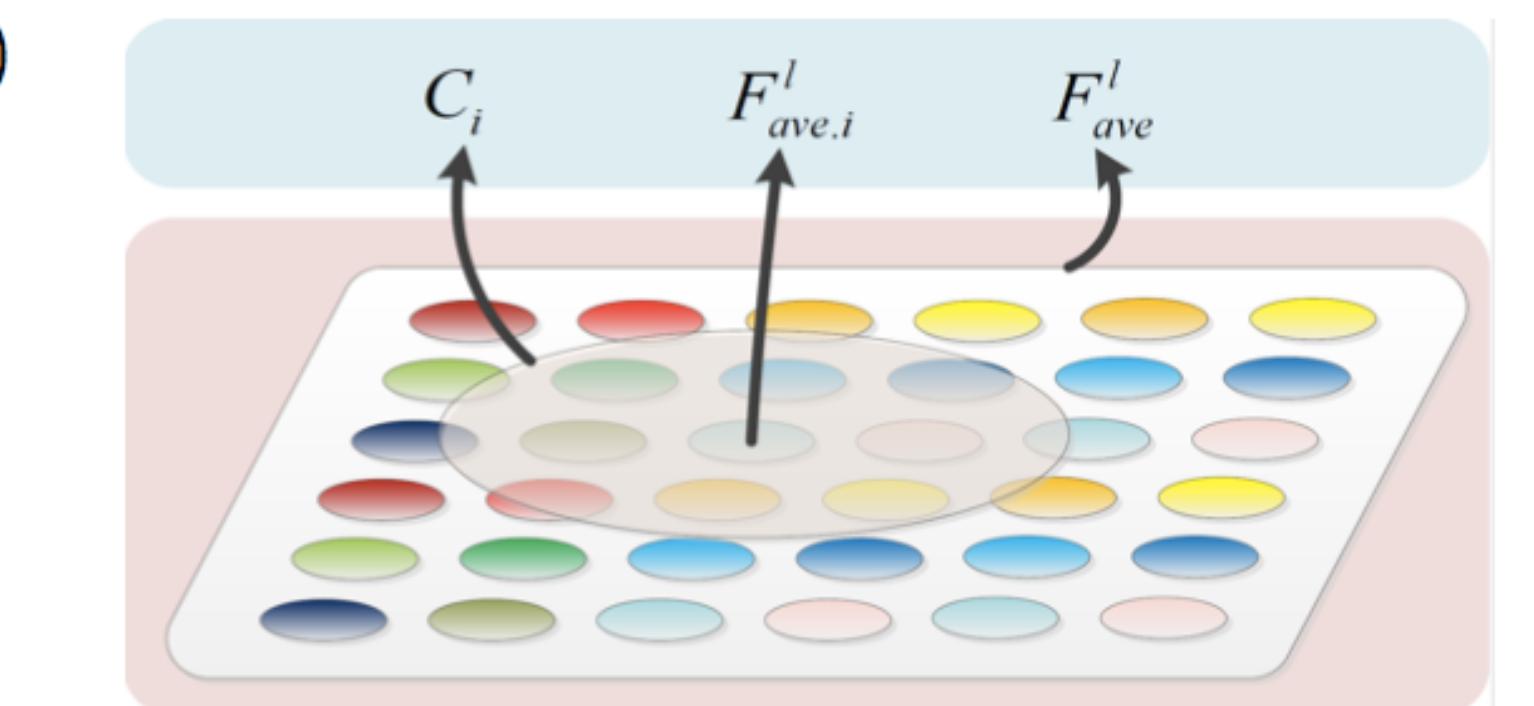


Figure 3. Estimation of  $\gamma_i$  and  $\delta_i^2$  with MLE.

$$\hat{\gamma}_i = 1/|C_i| \sum_{i \in C_i} F_{ave,i}^l,$$

$$\hat{\delta}_i^2 = 1/|C_i| \sum_{i \in C_i} (F_{ave,i}^l - \hat{\gamma}_i)(F_{ave,i}^l - \hat{\gamma}_i)^T$$

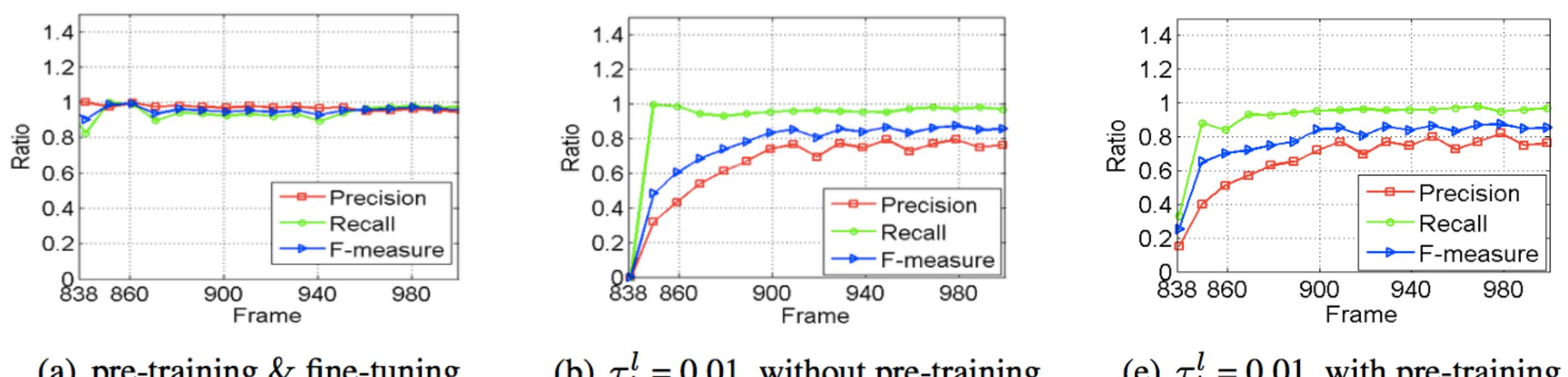
$$\tau_i^l = (\tau_{temporal,i}^l + \tau_{spatial,i}^l)/2$$

$$\hat{d}_i^l(t) = \min \{ \hat{d}_{i1}^l(t), \hat{d}_{i2}^l(t), \dots, \hat{d}_{iQ}^l(t) \}$$

## Evaluation on the CDnet 2014 Dataset :

Method	$FM_{overall}$	$FM_{BL}$	$FM_{CJ}$	$FM_{DB}$	$FM_{IOM}$	$FM_{SH}$	$FM_{TH}$	$FM_{BW}$	$FM_{LF}$	$FM_{NV}$	$FM_{PTZ}$	$FM_{TU}$
STSOM	<b>0.816</b>	<b>0.957</b>	<b>0.888</b>	<b>0.923</b>	<b>0.835</b>	<b>0.910</b>	<b>0.848</b>	<b>0.892</b>	<b>0.812</b>	<b>0.563</b>	<b>0.575</b>	<b>0.800</b>
IUTIS-5	<b>0.771</b>	<b>0.956</b>	<b>0.833</b>	0.890	0.729	<b>0.908</b>	0.830	0.824	<b>0.774</b>	0.529	0.428	<b>0.783</b>
SharedModel	0.747	0.952	0.814	0.822	0.672	0.845	0.831	0.798	0.728	0.541	0.386	0.733
SuBSENSE	0.741	0.950	0.815	0.817	0.656	0.864	0.817	<b>0.861</b>	0.644	0.559	0.347	0.779
PAWCS	0.740	0.939	0.813	<b>0.893</b>	0.776	0.871	0.832	0.815	0.658	0.415	0.461	0.645
C-EFIC	0.730	0.930	0.824	0.562	0.622	0.845	<b>0.834</b>	0.786	0.680	<b>0.667</b>	<b>0.620</b>	0.627
MBS	0.728	0.928	0.836	0.791	0.756	0.826	0.819	0.798	0.635	0.515	0.552	0.585
FTSG	0.728	0.933	0.751	0.879	<b>0.789</b>	0.853	0.776	0.822	0.625	0.513	0.324	0.712
S-Subsense	0.717	0.948	0.807	0.815	0.601	0.865	0.685	0.859	0.651	0.534	0.339	0.751
SMSOM	-	0.944	0.732	0.675	-	-	0.793	-	-	-	-	-
SOBS	0.596	0.933	0.705	0.643	0.562	0.721	0.683	0.662	0.546	0.450	0.040	0.488
KDE	0.568	0.909	0.572	0.596	0.408	0.766	0.742	0.757	0.547	0.436	0.036	0.447
GMM	0.556	0.838	0.596	0.633	0.520	0.715	0.662	0.738	0.537	0.409	0.152	0.466

## Evaluation of Pre-training and Fine-tuning :



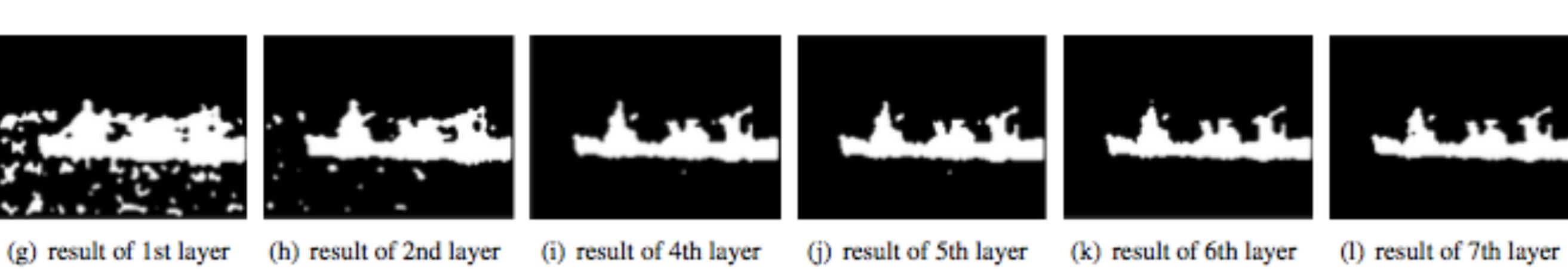
(a) pre-training & fine-tuning      (b)  $\tau_i^l = 0.01$ , without pre-training      (e)  $\tau_i^l = 0.01$ , with pre-training

## Evaluation of SSOM, TSOM and STSOM :



(a) original image      (b) ground truth      (c) with pre-training      (e) TSOM result      (f) SSOM result

## Evaluation of Deep Network :



(g) result of 1st layer      (h) result of 2nd layer      (i) result of 4th layer      (j) result of 5th layer      (k) result of 6th layer      (l) result of 7th layer