

UberNet

Kamil Bladoszewski

UberNet: Training a 'Universal' Convolutional Neural Network for Low-, Mid-, and High-Level Vision using Diverse Datasets and Limited Memory

Plan of presentation

1. Introduction
2. Tasks
3. Architecture
4. Multi-Task Training
5. Low memory usage
6. Results
7. Bibliography

Introduction

The *swiss knife* of Convolutional Neural Networks

Abstract

In this work we introduce a convolutional neural network (CNN) that jointly handles low-, mid-, and high-level vision tasks in a unified architecture that is trained end-to-end. Such a universal network can act like a ‘swiss knife’ for vision tasks; we call this architecture an UberNet to indicate its overarching nature.

The *swiss knife* of Convolutional Neural Networks



The *swiss knife* of Convolutional Neural Networks



Source: <https://oddtymall.com/includes/content/giant-swiss-army-knife-thumb.jpg>

The *swiss knife* of Convolutional Neural Networks

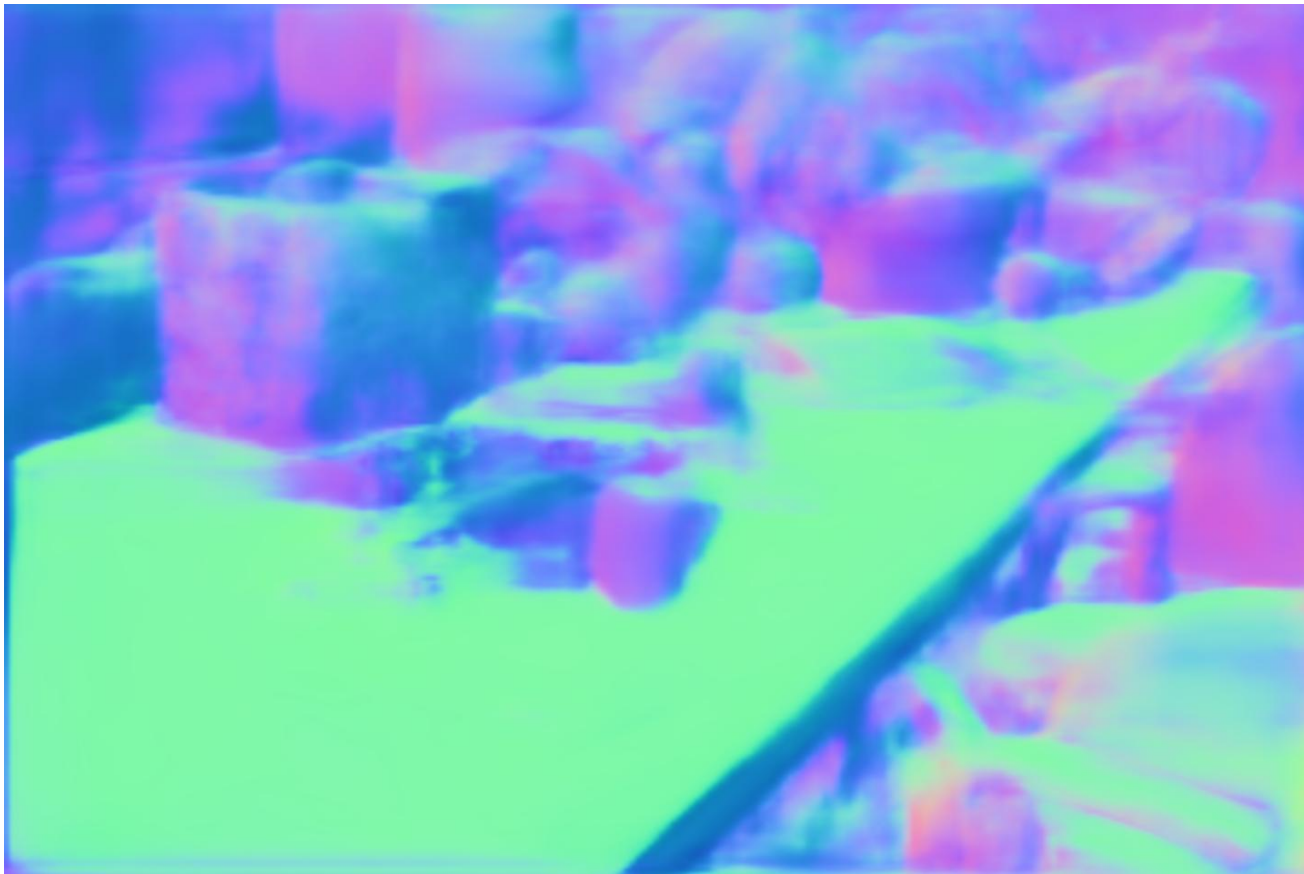


Tasks

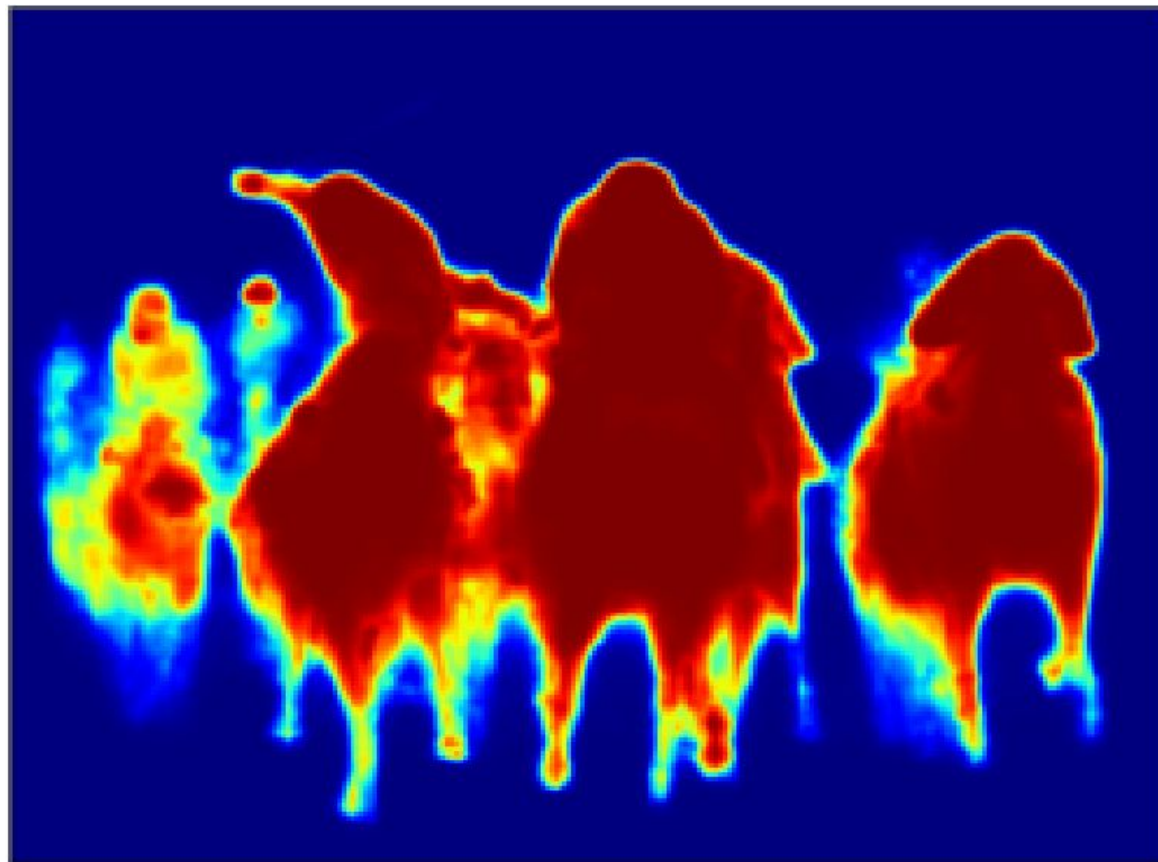
Boundary Detection



Normal Estimation



Saliency Estimation



Source: *UberNet: Training a 'Universal' Convolutional Neural Network for Low-, Mid-, and High-Level Vision using Diverse Datasets and Limited Memory* <https://arxiv.org/abs/1609.02132>

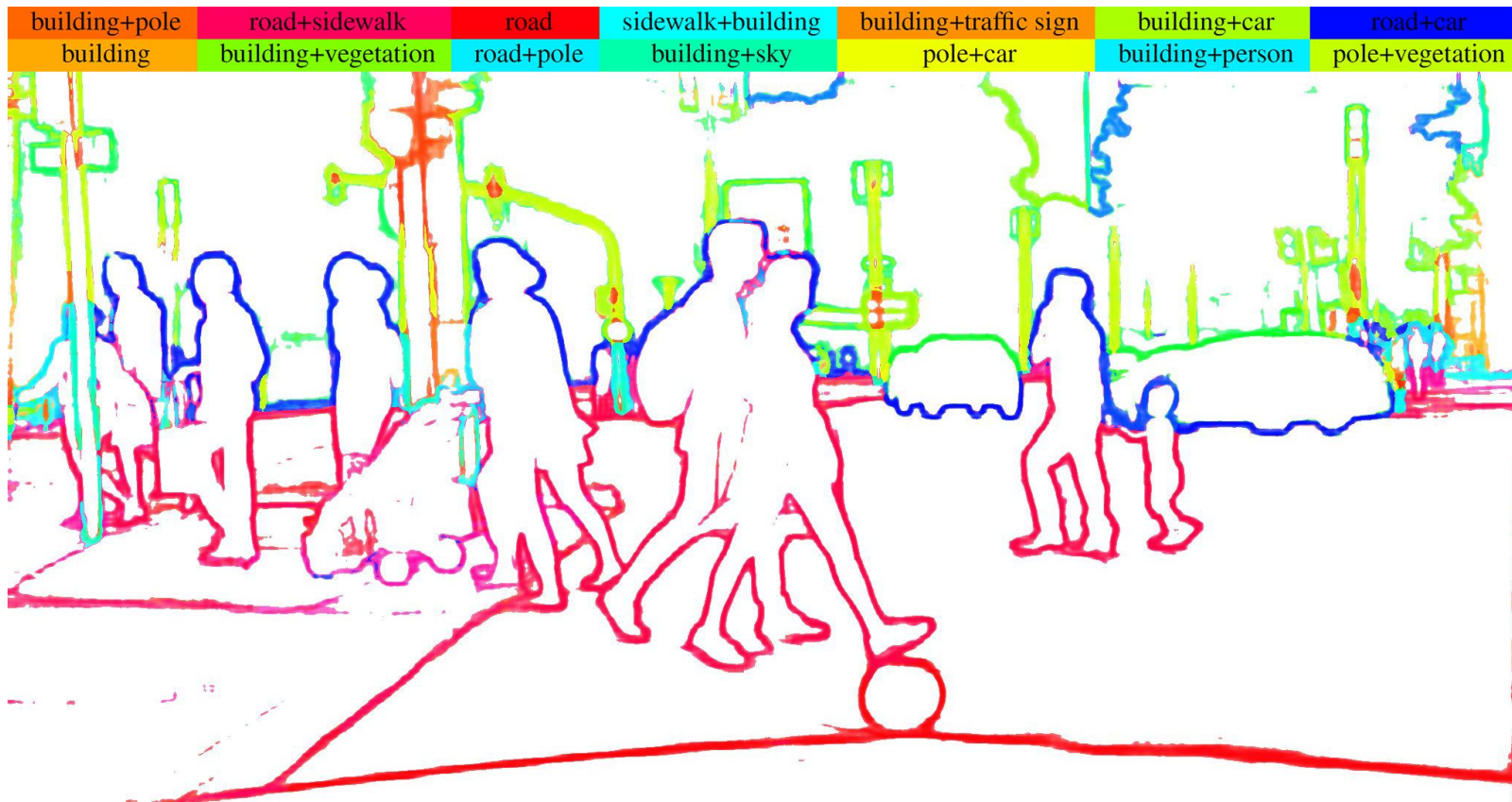
Semantic Segmentation



Human Part Segmentation



Semantic Boundary Detection

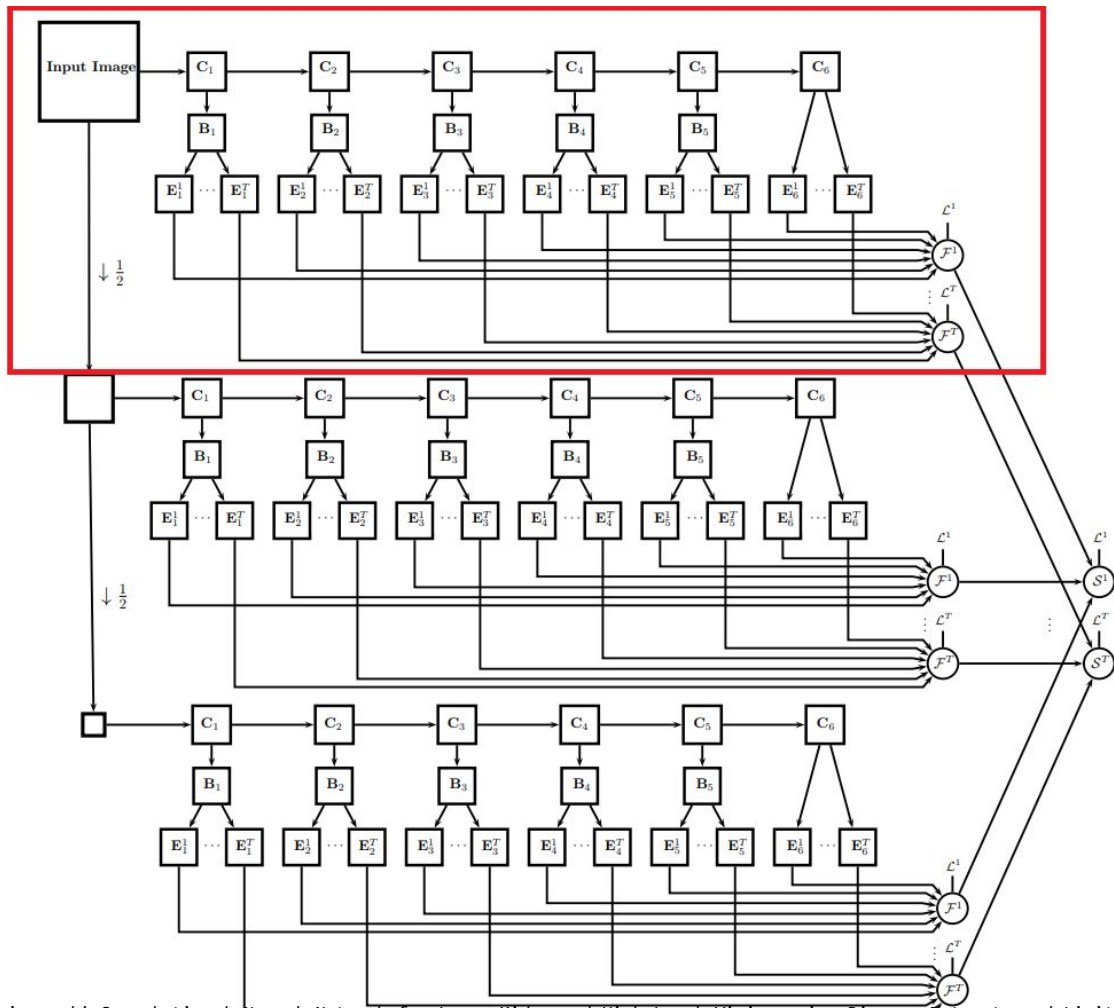


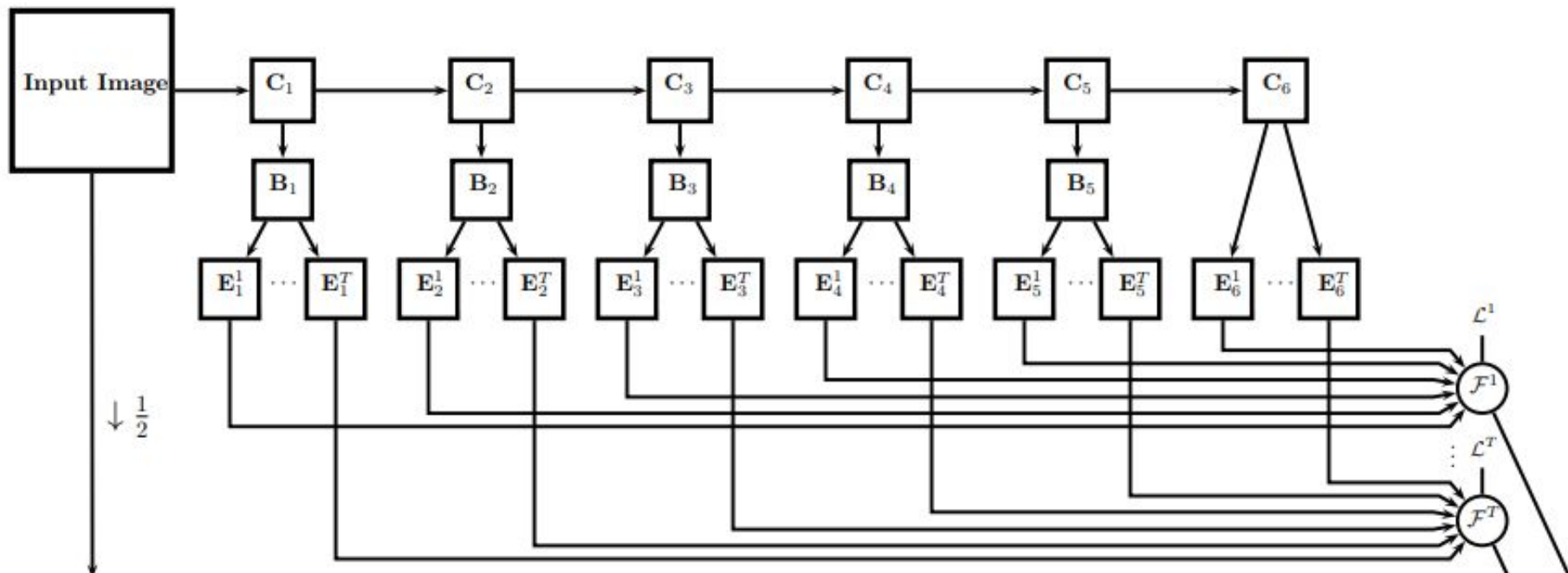
Region Proposal Generation and Object Detection



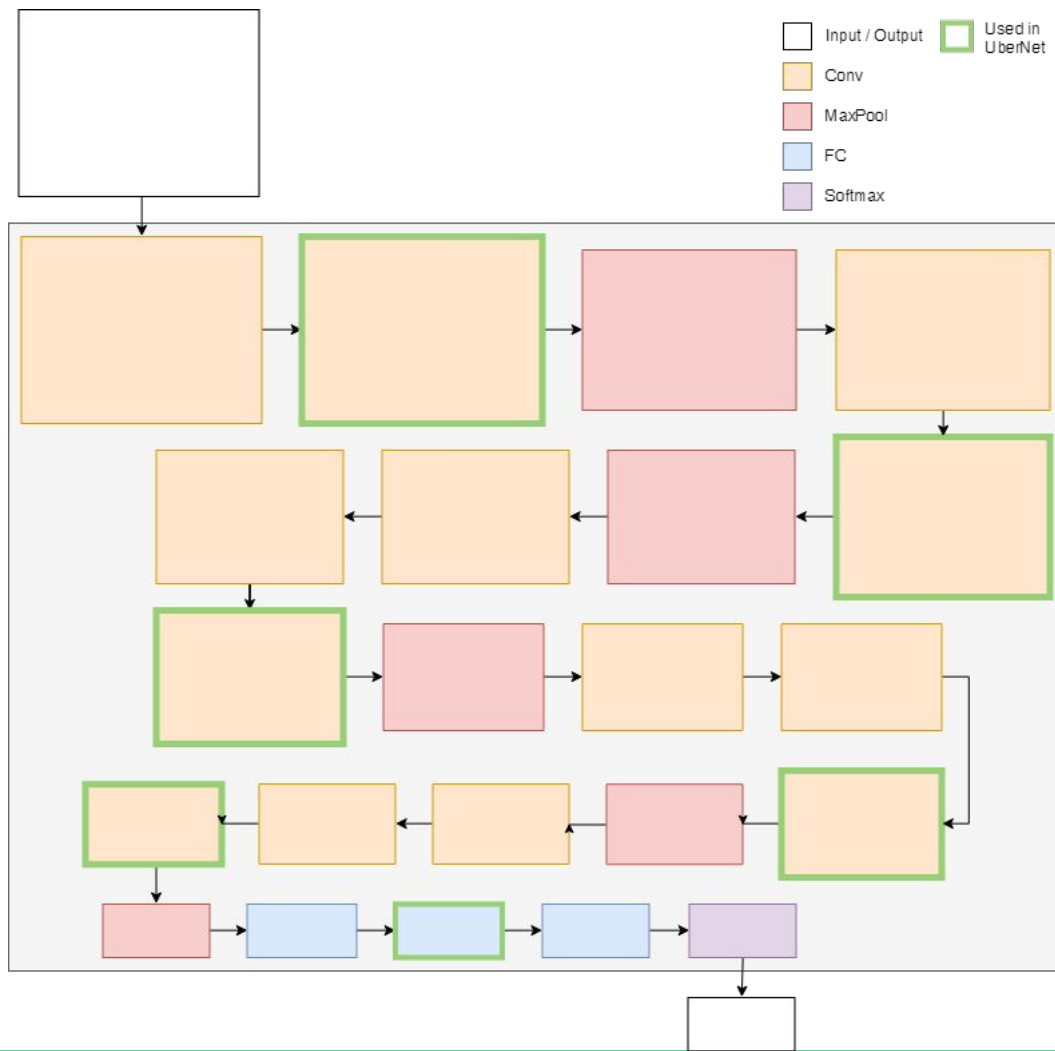
Architecture

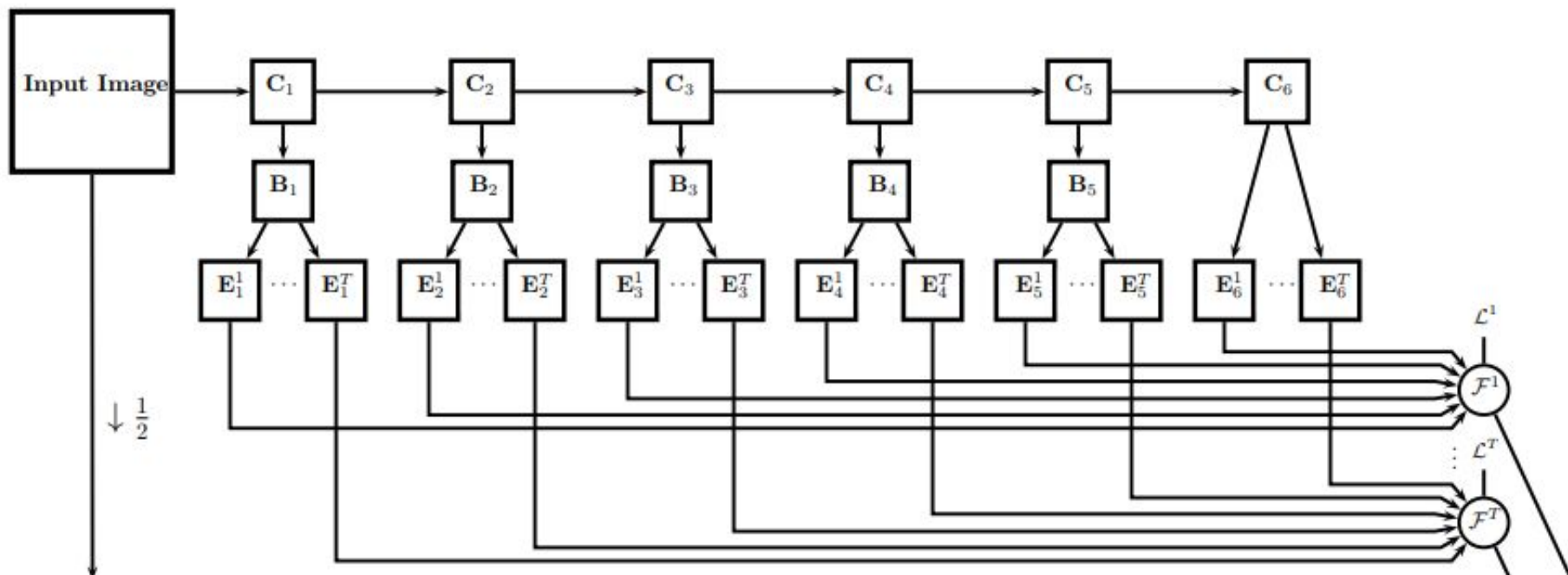


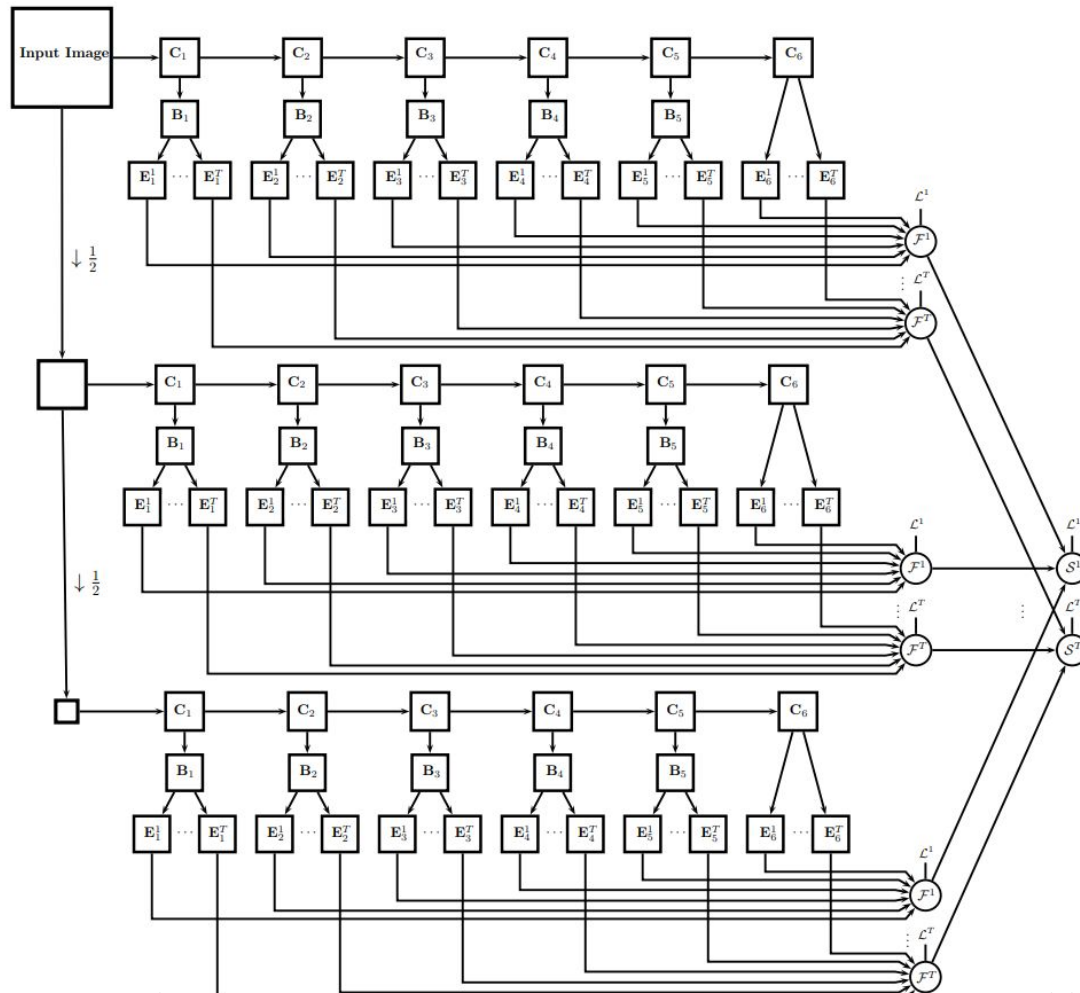


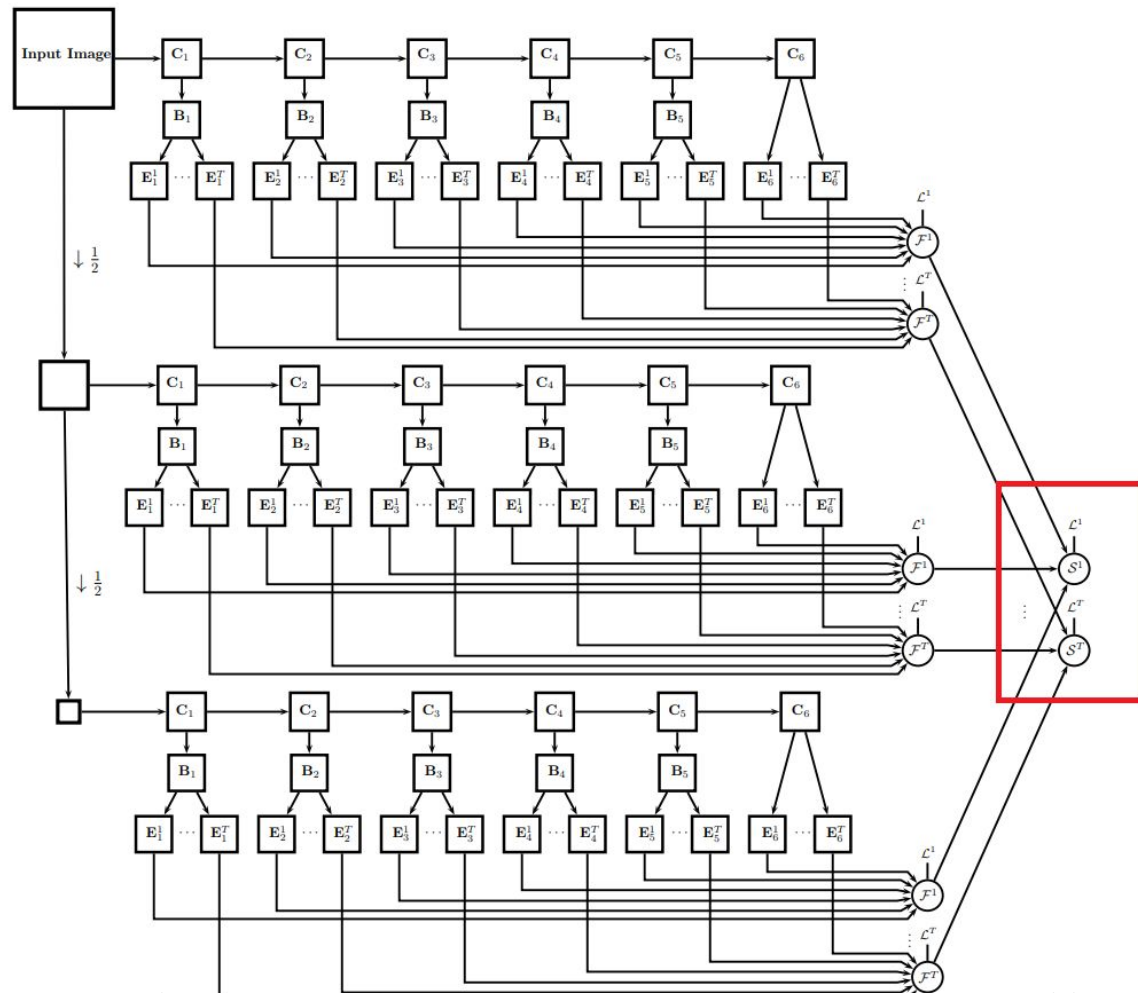


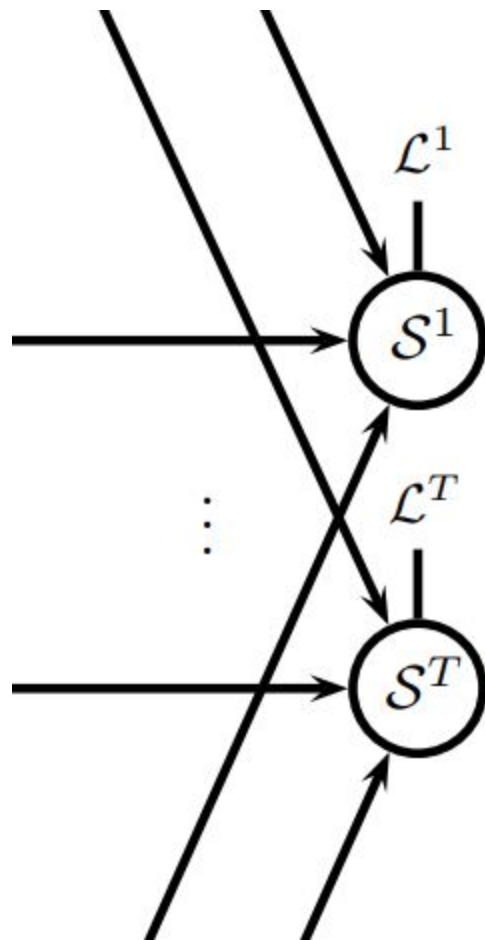
VGG-16

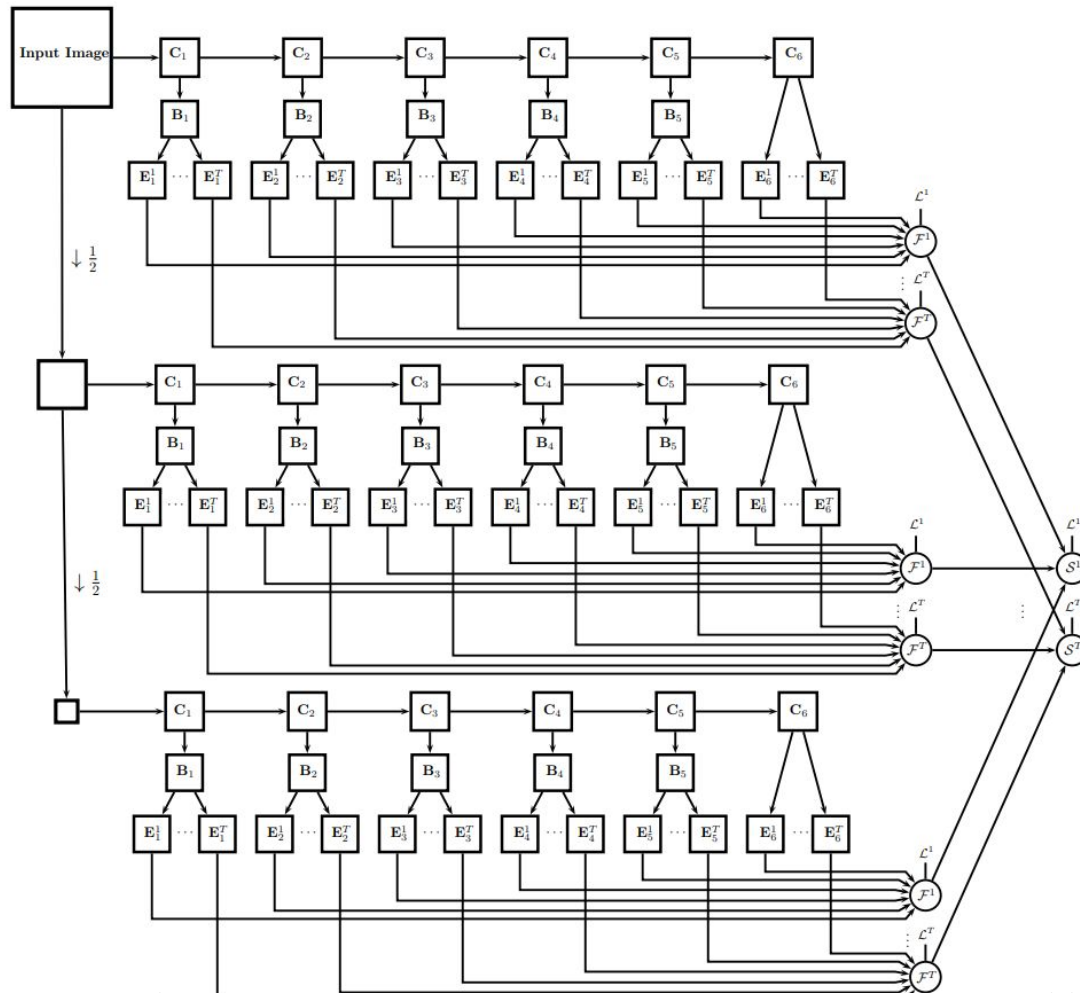




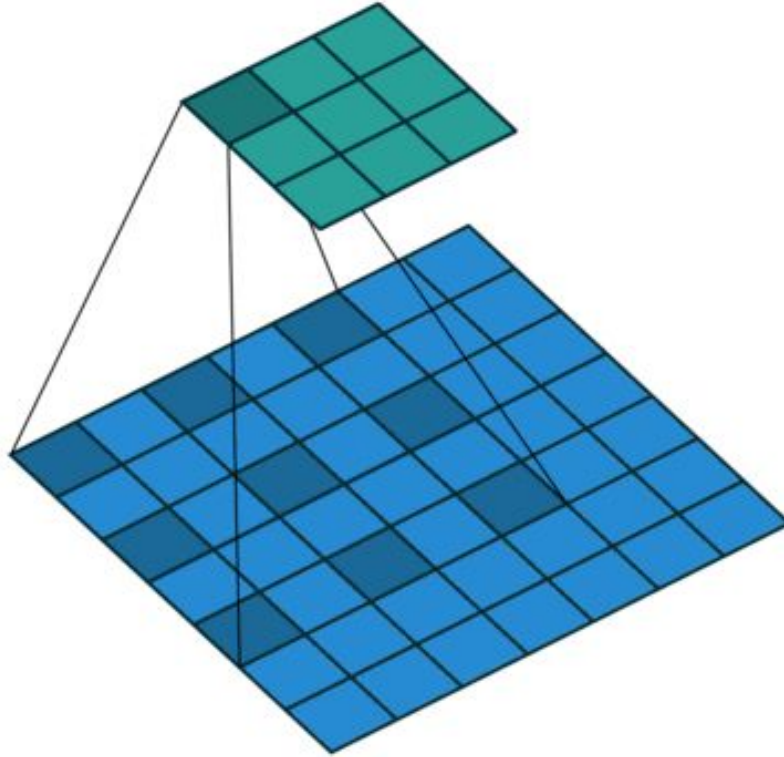








Atrous/Dilated Convolutions



Architecture

- Skip layers
- Skip-layer normalization
- Cumulative task-specific operations
- Fusion layers
- Atrous convolution
- Multi-resolution CNN
- Task-specific deviations

Multi-Task Training

Objective: minimize the loss

$$\mathcal{L}(\mathbf{w}_{0,t_1,\dots,t_T}) = \mathcal{R}(\mathbf{w}_0) + \sum_{t=1}^T \gamma_{t_k} (\mathcal{R}(\mathbf{w}_t) + L_t(\mathbf{w}_0, \mathbf{w}_t))$$

$$L_t(\mathbf{w}_0, \mathbf{w}_t) = \frac{1}{N} \sum_{i=1}^N \delta_{t,i} L_t(\mathbf{f}_t^i(\mathbf{w}_0, \mathbf{w}_t), \mathbf{y}_t^i)$$

Many task = many different type annotations

Problem: Diverse datasets

There isn't a dataset,
which contains annotations for all the problems

Divers datasets

	Imagenet [88]	VOC'07 [26]	VOC'10 [26]	VOC'12 [26]	MS-COCO [62]	NYU [74]	MSRA10K [18]	BSD [69]
Detection	Partial	Yes	Yes	Yes	Yes	No	No	No
Semantic Segmentation	No	Partial	[36, 71]	Partial	Yes	Yes*	No	No
Instance Segmentation	No	Partial	[36, 71]	Partial	Yes	No	No	No
Human parts	No	No	[17]	No	No	No	No	No
Human landmarks	No	No	[10]	No	Yes	No	No	No
Surface Normals	No	No	No	No	No	Yes	No	No
Saliency	No	No	No	No	No	No	Yes	No
Boundaries	No	No	[71]	No	No	No	No	Yes
Symmetry	No	No	No	No	Partial, [91]	No	No	[97]

Divers datasets

	VOC'07	VOC'12	VOC'12	NYU	MSRA10K	BSD
	trainval	train	val			
	5011	5717	5823	23024	10000	5100
Detection	5011	5717	5823	0	0	0
S. Segmentation	422	4998	5105	0	0	0
S. Boundaries	0	4998	5105	0	0	0
Human Parts	0	1716	1817	0	0	0
Normals	0	0	0	23024	0	0
Saliency	0	0	0	0	10000	0
Boundaries	0	4998	5105	0	0	5100

PASCAL in Detail

	10100
Semantic Segmentation	10100
Human Parts	3548
Occlusion	10094
Boundaries	10100

Many task = many different annotation types

Problem: Diverse datasets

There isn't a dataset,
which contains annotations for all the problems

Synchronous SGD

```
for  $m = 1$  to  $M$  do  
  {construct minibatch}  
   $\mathcal{B} \leftarrow \{i_1, \dots, i_B\}$  with  $i_i \sim U[1, N]$   
  {initialize gradient accumulators}  
   $\mathbf{dw}_0 \leftarrow 0, \mathbf{dw}_1 \leftarrow 0, \dots, \mathbf{dw}_T \leftarrow 0$   
  for  $i \in \mathcal{B}$  do  
    {cnn gradients}  
     $\mathbf{dw}_0 \leftarrow \mathbf{dw}_0 + \sum_t \delta_{t,i} \gamma_t \nabla_{\mathbf{w}_0} L_t (\mathbf{f}_t^i(\mathbf{w}_0, \mathbf{w}_t), \mathbf{y}_t^i)$   
    {Task gradients,  $t = 1, \dots, T$ }  
     $\mathbf{dw}_t \leftarrow \mathbf{dw}_t + \delta_{t,i} \gamma_t \nabla_{\mathbf{w}_0} L_t (\mathbf{f}_t^i(\mathbf{w}_0, \mathbf{w}_t), \mathbf{y}_t^i)$   
  end for  
  for  $p \in \{0, 1, \dots, T\}$  do  
     $\mathbf{w}_p \leftarrow \mathbf{w}_p - \epsilon (\lambda \mathbf{w}_p + \frac{1}{B} \mathbf{dw}_p)$   
  end for  
end for
```

Asynchronous SGD

```

{initialize gradient accumulators}
 $\mathbf{dw}_0 \leftarrow 0, \mathbf{dw}_1 \leftarrow 0, \dots, \mathbf{dw}_T \leftarrow 0$ 
{initialize counters}
 $\mathbf{c}_0 \leftarrow 0, \mathbf{c}_1 \leftarrow 0, \dots, \mathbf{c}_T \leftarrow 0$ 
for  $m = 1$  to  $M \cdot B$  do
    Sample  $i \sim U[1, N]$ 
    {cnn gradients & counter: always updated}
     $\mathbf{c}_0 \leftarrow \mathbf{c}_0 + 1$ 
     $\mathbf{dw}_0 \leftarrow \mathbf{dw}_0 + \sum_t \delta_{t,i} \gamma_t \nabla_{\mathbf{w}_0} L_t(\mathbf{f}_t^i(\mathbf{w}_0, \mathbf{w}_t), \mathbf{y}_t^i)$ 
    for  $t \in \{1, \dots, T\}$  do
        if  $\delta_{t,i} = 1$  then
            {update accumulator and counter for task  $t$  if the
             current sample is relevant}
             $\mathbf{c}_t \leftarrow \mathbf{c}_t + 1$ 
             $\mathbf{dw}_t \leftarrow \mathbf{dw}_t + \gamma_t \nabla_{\mathbf{w}_t} L_t(\mathbf{f}_t^i(\mathbf{w}_0, \mathbf{w}_t), \mathbf{y}_t^i)$ 
        end if
    end for
end for

```

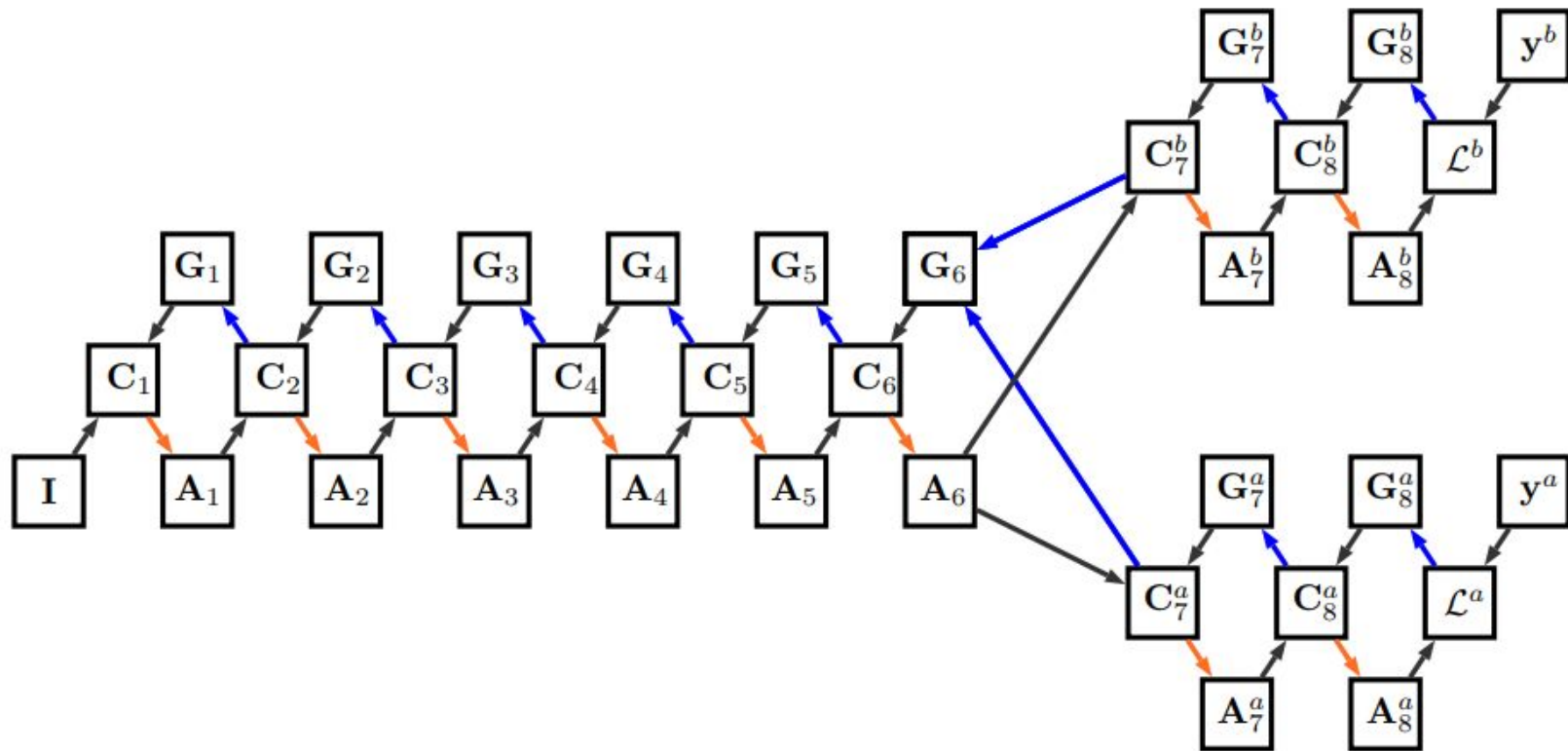
```

for  $p \in \{0, 1, \dots, T\}$  do
    if  $\mathbf{c}_p = B_p$  then
        {update parameters if we have seen enough}
         $\mathbf{w}_p \leftarrow \mathbf{w}_p - \epsilon \left( \lambda \mathbf{w}_p + \frac{1}{B_p} \mathbf{dw}_p \right)$ 
         $\mathbf{c}_p \leftarrow 0, \mathbf{dw}_p \leftarrow 0$ 
    end if
end for

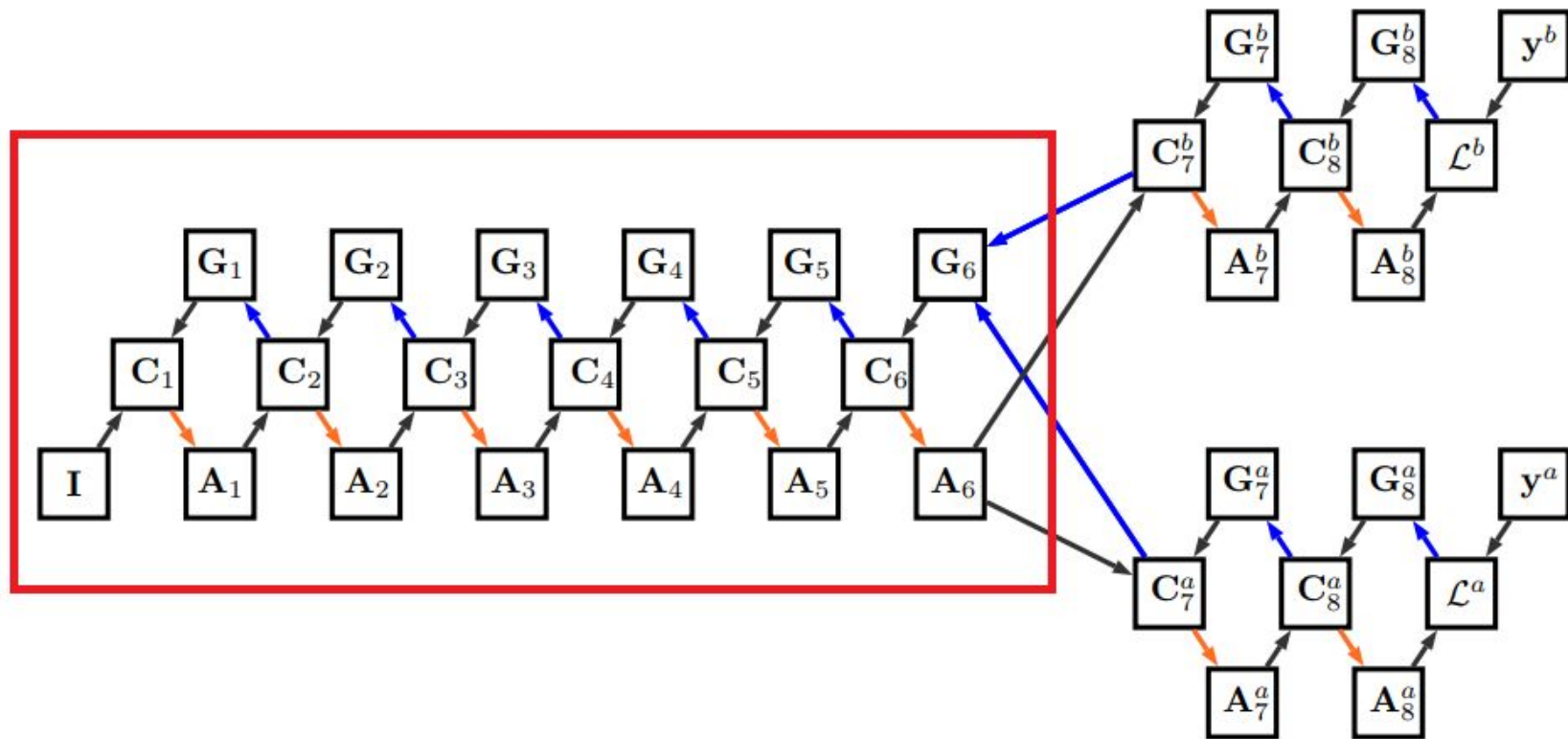
```

Low memory usage

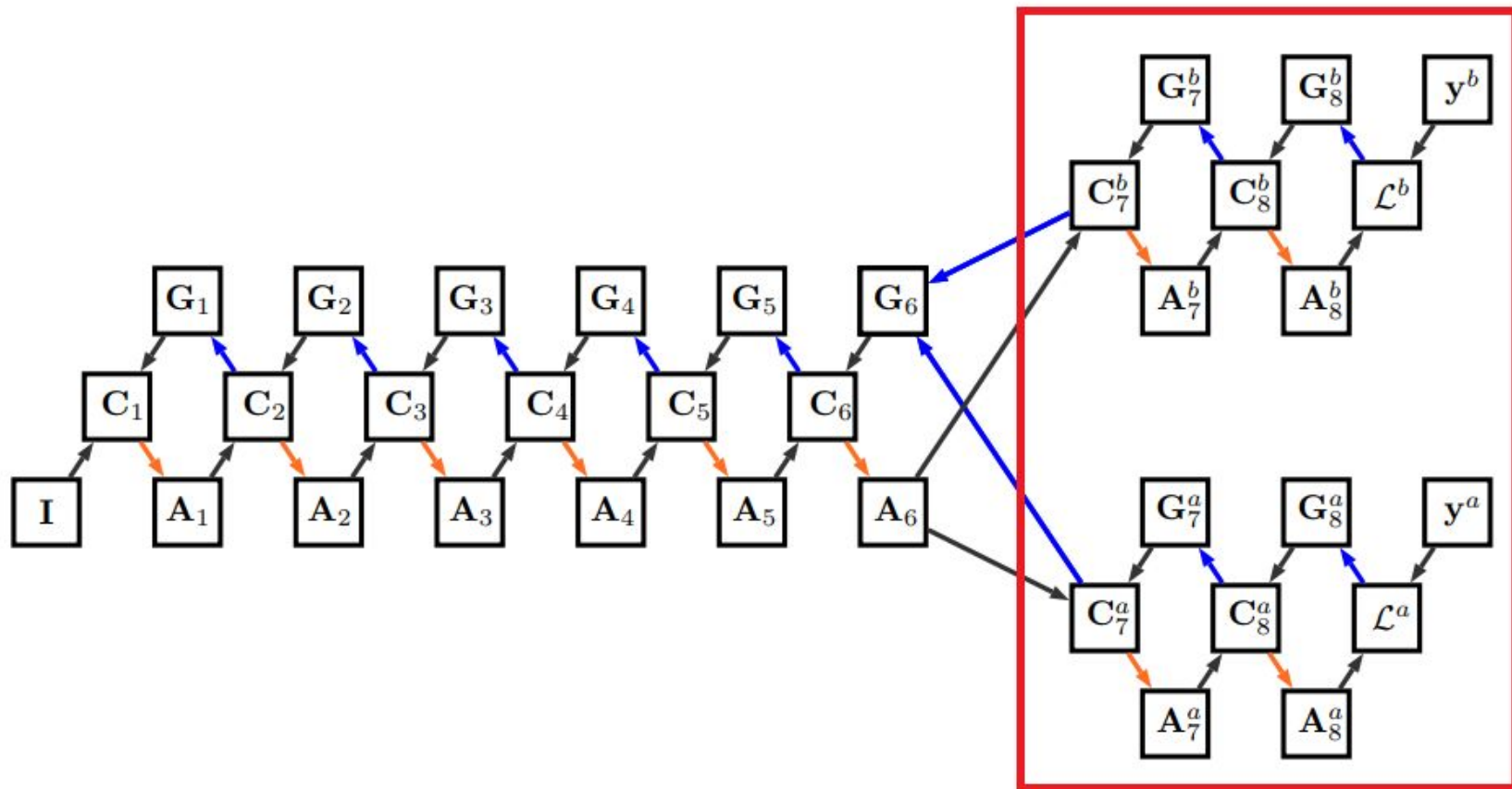
Normal memory forward and backward pass



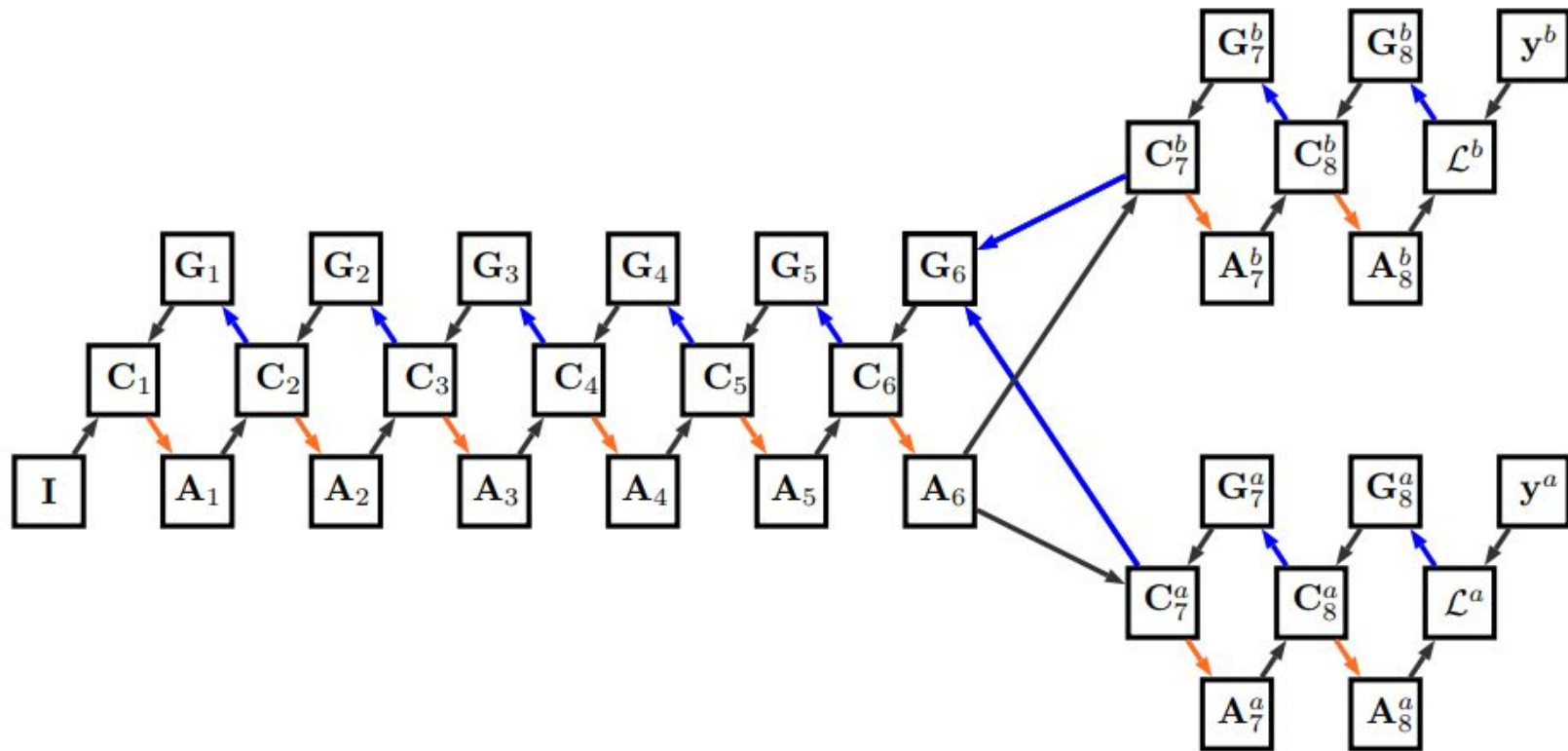
Normal memory forward and backward pass



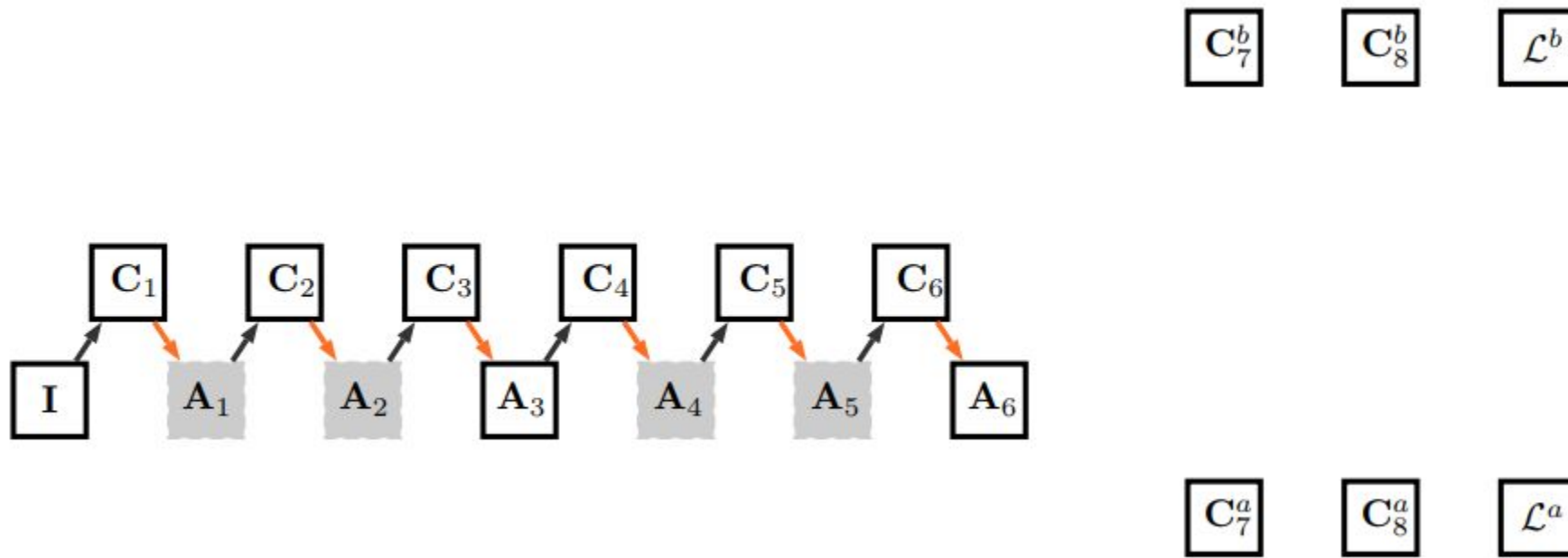
Normal memory forward and backward pass



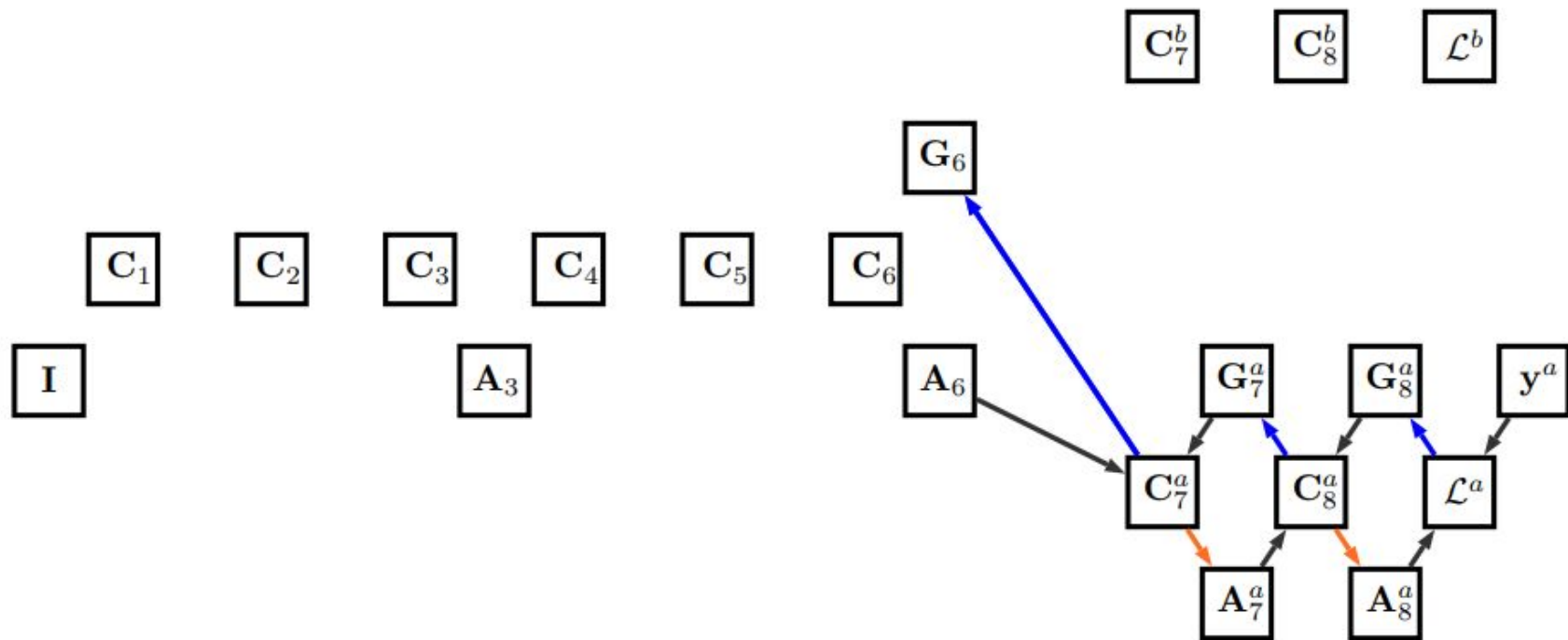
Normal memory forward and backward pass



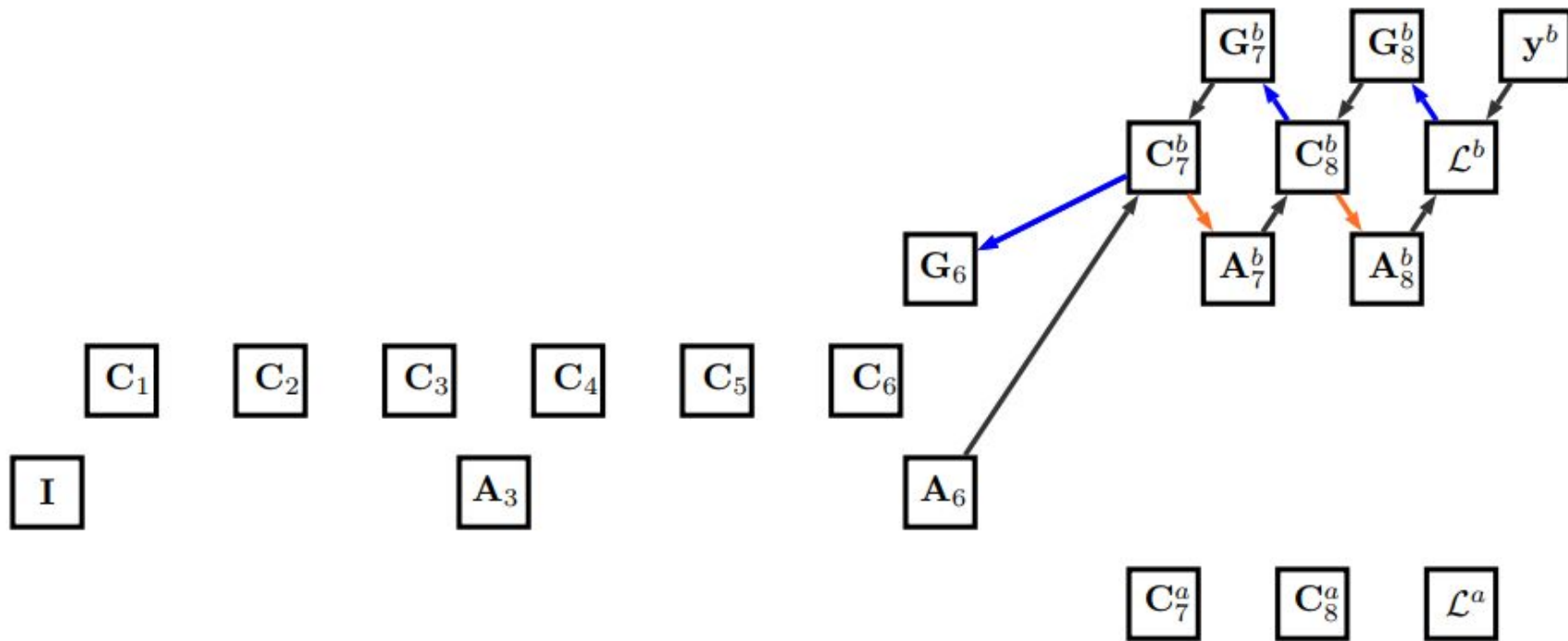
Low memory forward and backward pass



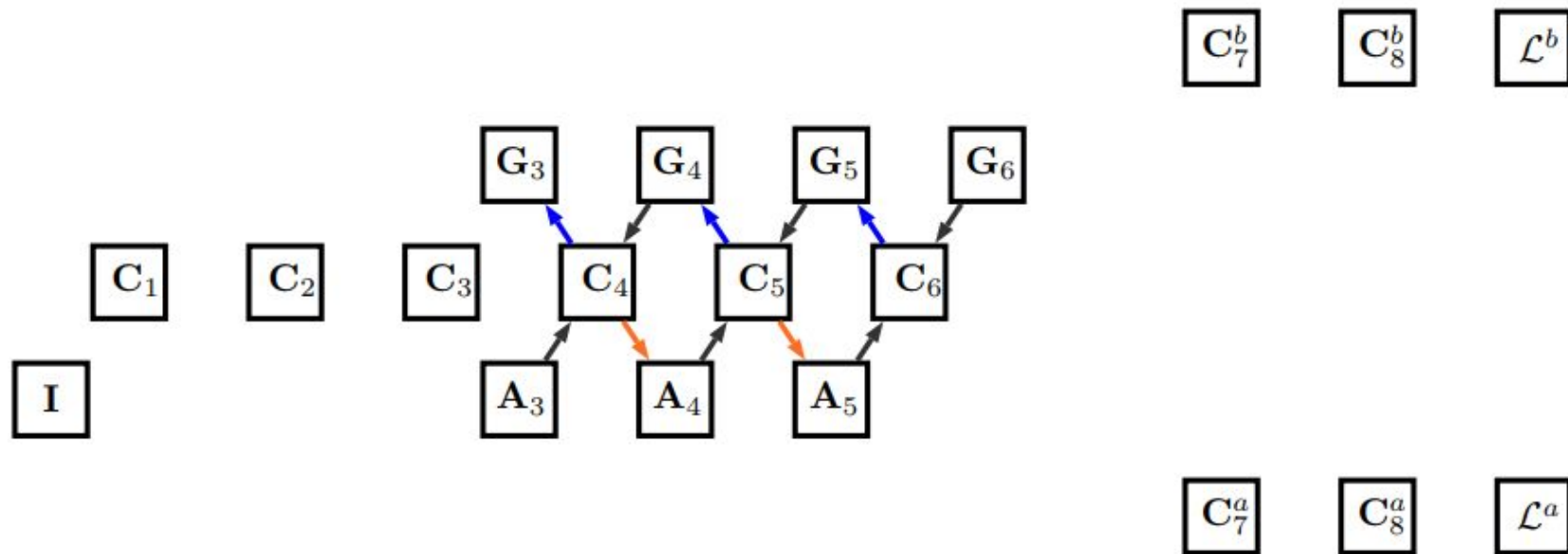
Low memory forward and backward pass



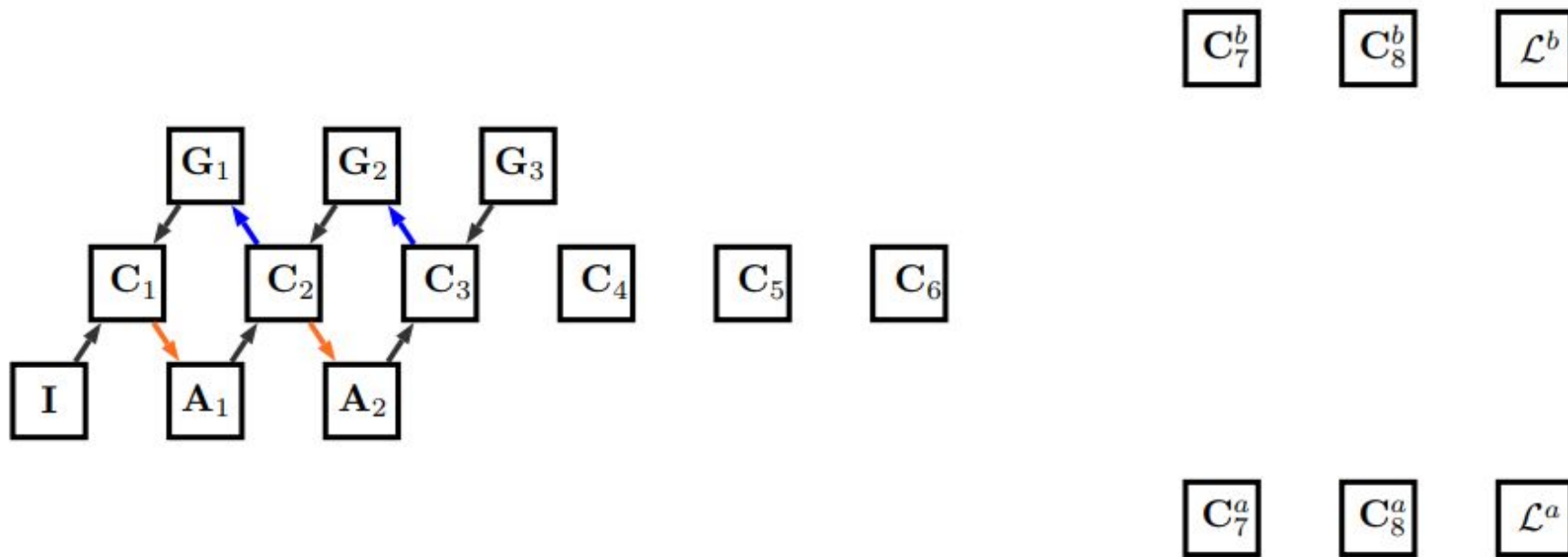
Low memory forward and backward pass



Low memory forward and backward pass



Low memory forward and backward pass



Results

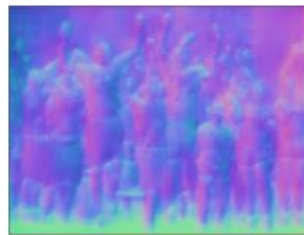
Input



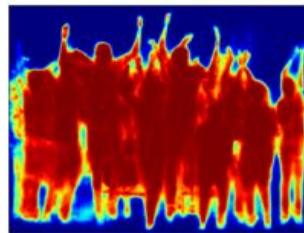
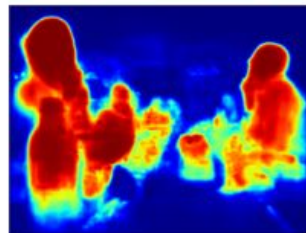
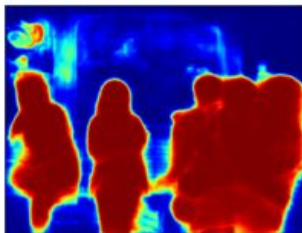
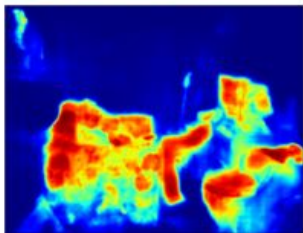
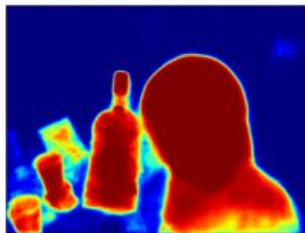
Boundaries



Surface Normals



Saliency

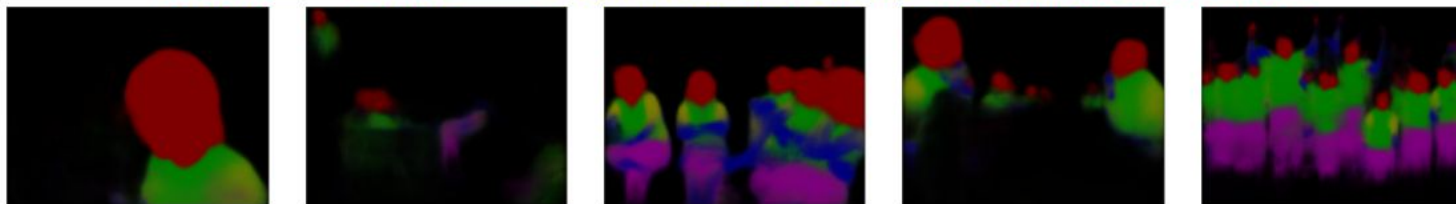


Object Detection Sem.c Segmentation Sem.c Boundaries

airplane	bicycle	bird	boat	bottle	bus	car	cat	chair	cow
d.table	dog	horse	m.bike	person	p.plant	sheep	sofa	train	tv



head	torso	upper arm	lower arm	upper leg	lower leg
------	-------	-----------	-----------	-----------	-----------

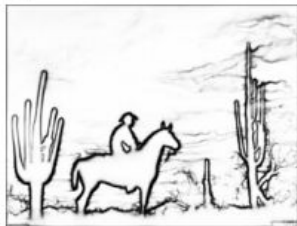


Human Parts

Input



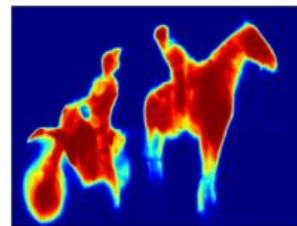
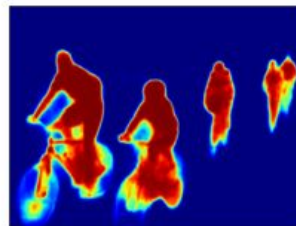
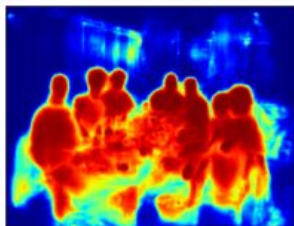
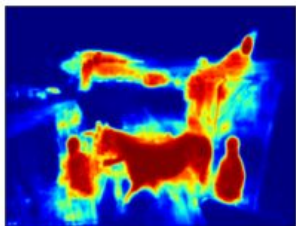
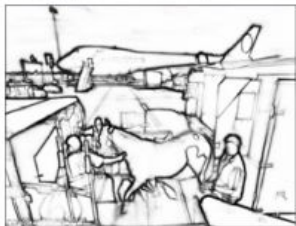
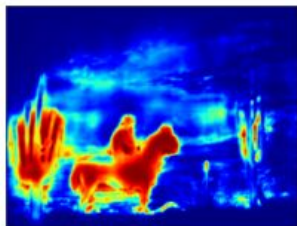
Boundaries



Surface Normals

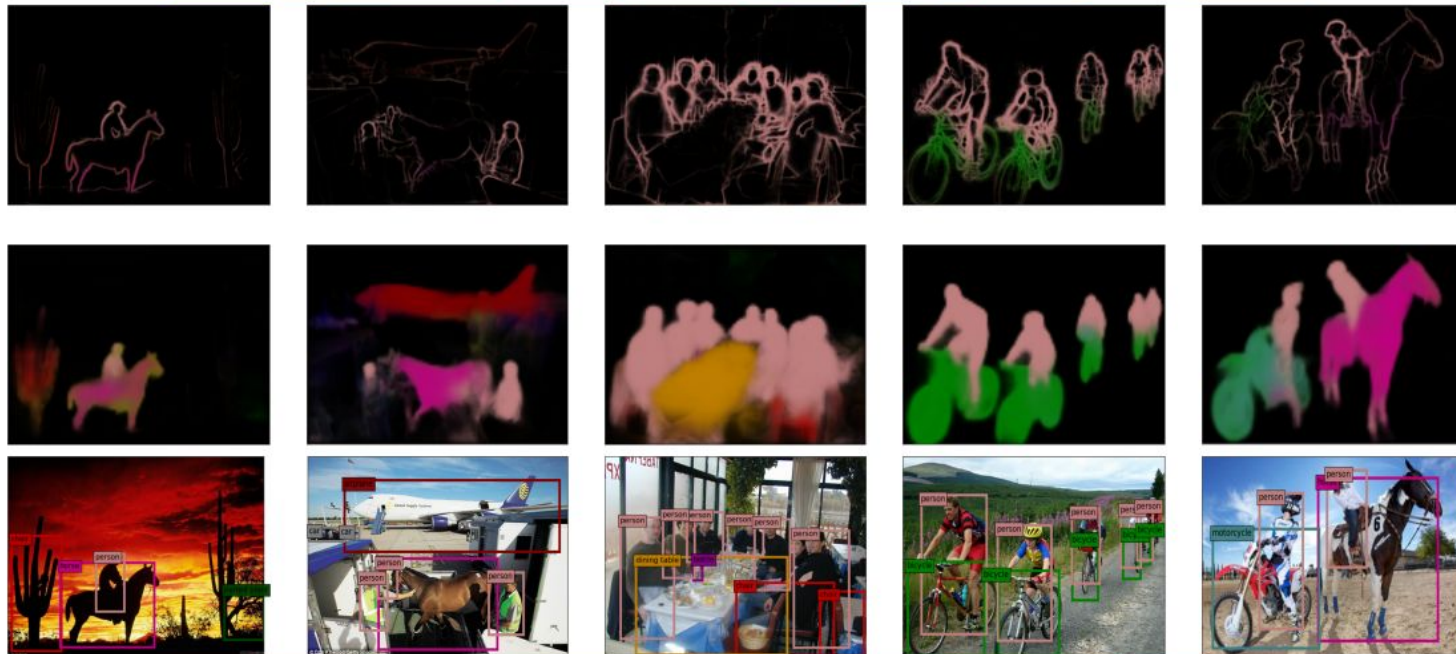


Saliency

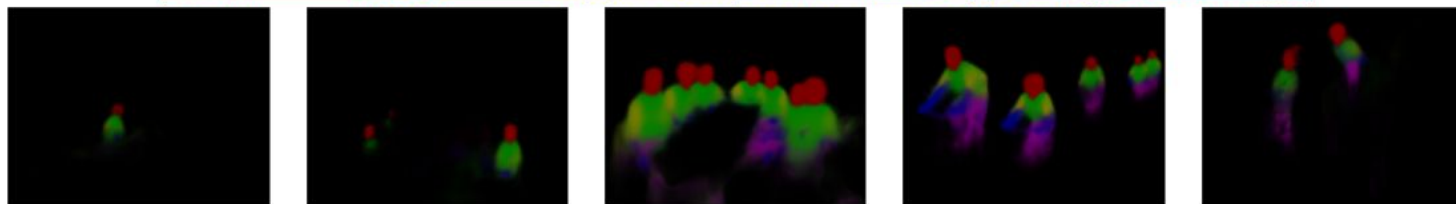


Object Detection Sem.c Segmentation Sem.c Boundaries

airplane	bicycle	bird	boat	bottle	bus	car	cat	chair	cow
d.table	dog	horse	m.bike	person	p.plant	sheep	sofa	train	tv



head	torso	upper arm	lower arm	upper leg	lower leg
------	-------	-----------	-----------	-----------	-----------



Bibliography

- UberNet: Training a ‘Universal’ Convolutional Neural Network for Low-, Mid-, and High-Level Vision using Diverse Datasets and Limited Memory <https://arxiv.org/abs/1609.02132>
- Very Deep Convolutional Networks for Large-Scale Image Recognition <https://arxiv.org/abs/1409.1556>
- CVPR’17 PASCAL in Detail Challenge <https://sites.google.com/view/pasd>
- Wikipedia:
 - Saliency map [https://en.wikipedia.org/wiki/Normal_\(geometry\)](https://en.wikipedia.org/wiki/Normal_(geometry))
 - Normal (geometry) https://en.wikipedia.org/wiki/Saliency_map
- An Introduction to different Types of Convolutions in Deep Learning <https://towardsdatascience.com/types-of-convolutions-in-deep-learning-717013397f4d>

Thank you!
