



Soutenance de Stage de L3

Multimodal Learning:

A case study for Gesture and Audio-Visual Speech Recognition

Hsieh Yu-Guan (Info 2016)

Supervised by
Amélie Cordier & Mathieu Lefort

Internship period: 14th June 2017 – 11th August 2017

behaviors.ai

Introduction

- BEHAVIORS.AI
- Hoomano & LIRIS
- TensorFlow



Introduction

- Artificial Intelligence

Introduction

- Artificial Intelligence
- > Robotics (embodied paradigm)

Introduction

- Artificial Intelligence
 - > Robotics (embodied paradigm)
 - > Developmental robotics

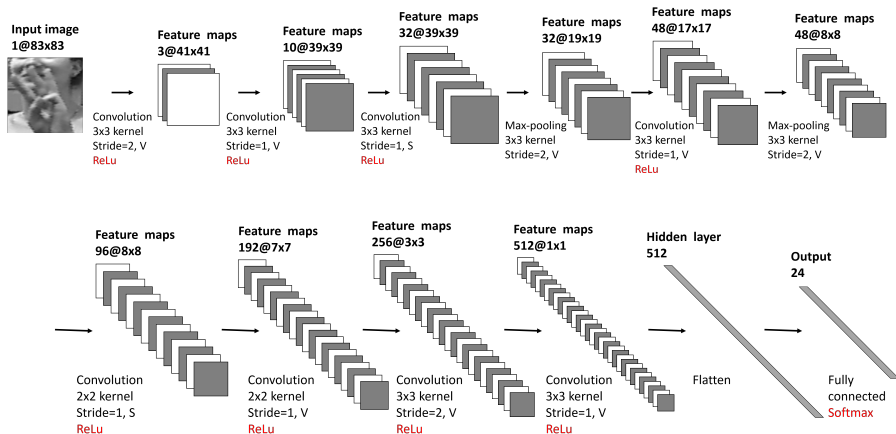
Introduction

- Artificial Intelligence
 - > Robotics (embodied paradigm)
 - > Developmental robotics
 - > Multimodal learning

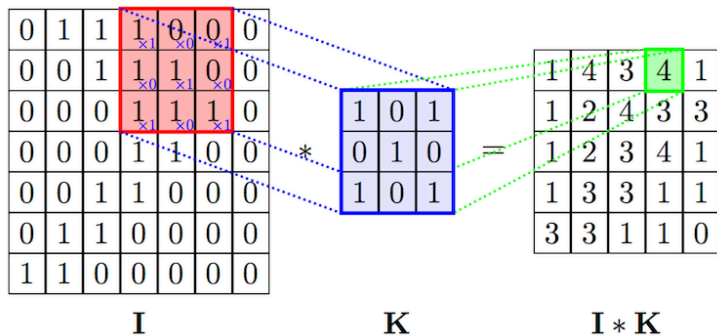
Introduction

- Artificial Intelligence
 - > Robotics (embodied paradigm)
 - > Developmental robotics
 - > Multimodal learning
 - > Gesture and Audio-Visual recognition

Deep Network Architectures – Convolutional Neural Networks (CNNs)



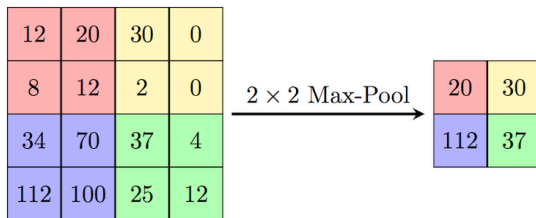
Deep Network Architectures – Convolutional Neural Networks (CNNs) – Convolution



Source: <https://cambridgespark.com/content/tutorials/convolutional-neural-networks-with-keras/index.html>

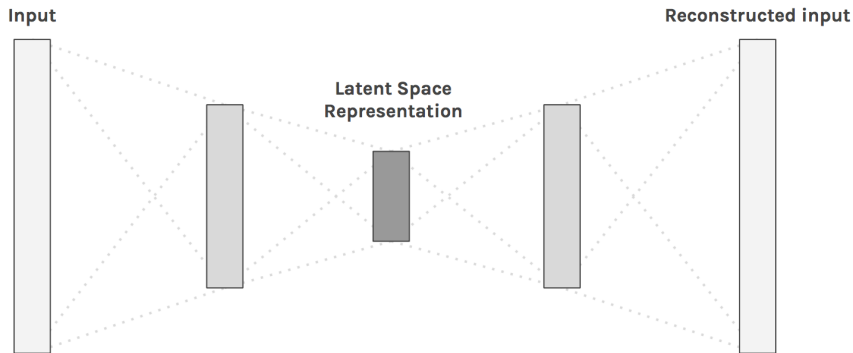
Also see https://github.com/vdumoulin/conv_arithmetic

Deep Network Architectures – Convolutional Neural Networks (CNNs) – Max-pooling



Source: <https://cambridgespark.com/content/tutorials/convolutional-neural-networks-with-keras/index.html>

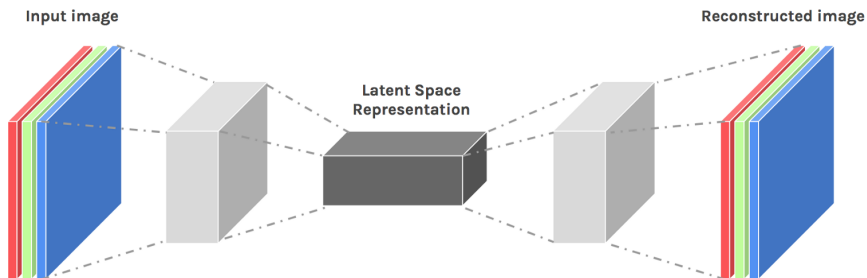
Deep Network Architectures – Autoencoder



Source:

<https://hackernoon.com/autoencoders-deep-learning-bits-1-11731e200694>

Deep Network Architectures – Convolutional Autoencoder



Source:

<https://hackernoon.com/autoencoders-deep-learning-bits-1-11731e200694>

Training a Machine Learning Model

- Loss function: cross-entropy, L2-distance
- Stochastic gradient descent (SGD)
- Backpropagation
- Variants of SGD: AdaGrad, Adam

Datasets – ASL Finger Spelling

- RGB and depth
- More than 60000 images for each modality
- 24 static signs in American Sign Language
- 5 subjects
- Only one channel in input
- Resized to 83×83 and Z-normalization



(a)



(b)

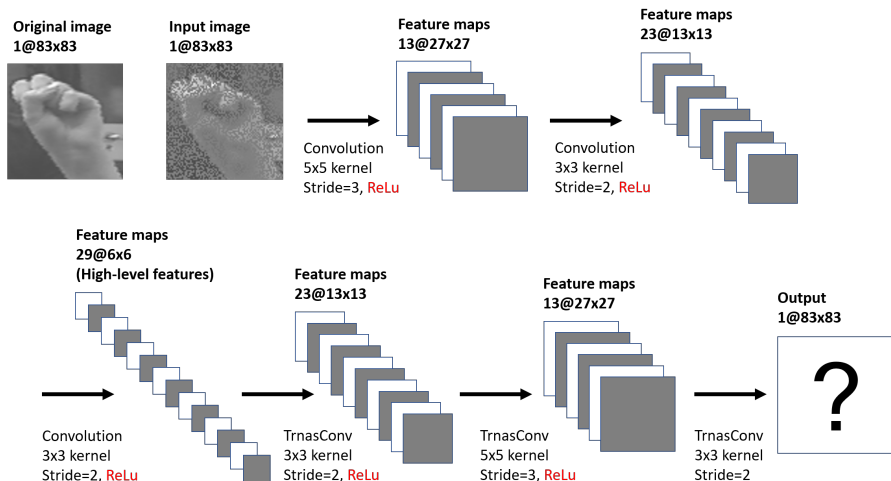


(c)

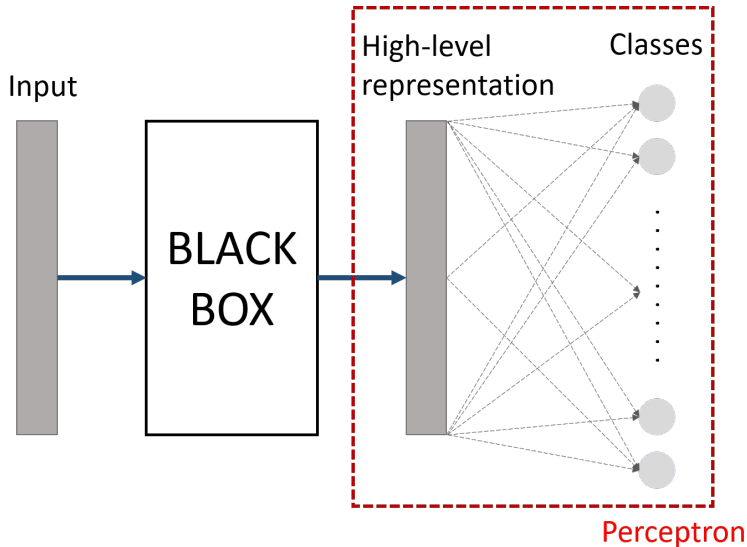


(d)

Results – Unsupervised Learning with CAE



Results – Unsupervised Learning with CAE



Results – Unsupervised Learning with CAE

- Raw: Perceptron that reads raw input data.
- CAE features: Perceptron stacked on the middle layer of the CAE.
- CAE architecture: Perceptron stacked on the middle layer of the CAE but train the whole network in a supervised way as a CNN.

		Raw	CAE features	CAE architecture
Intensity	train	69.47 %	78.87 %	91.29 %
	test	32.64 %	50.24 %	65.44 %
Depth	train	63.64 %	79.61 %	88.80 %
	test	29.93 %	41.64 %	55.62 %

Results – Shared Representation Learning

Input images

1@83x83

RGB(gray)



Depth



Feature maps
13@27x27



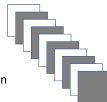
Feature maps
23@13x13



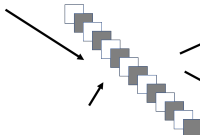
Convolution
5x5 kernel
Stride=3, ReLu



Convolution
3x3 kernel
Stride=2, ReLu

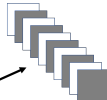


Feature maps
29@6x6
(Shared representation)

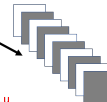


Convolution
3x3 kernel
Stride=2, ReLu

Feature maps
23@13x13



TrnasConv
3x3 kernel
Stride=2, ReLu



Feature maps
13@27x27



TrnasConv
3x3 kernel
Stride=2, ReLu



Outputs
1@83x83



TrnasConv
5x5 kernel
Stride=3

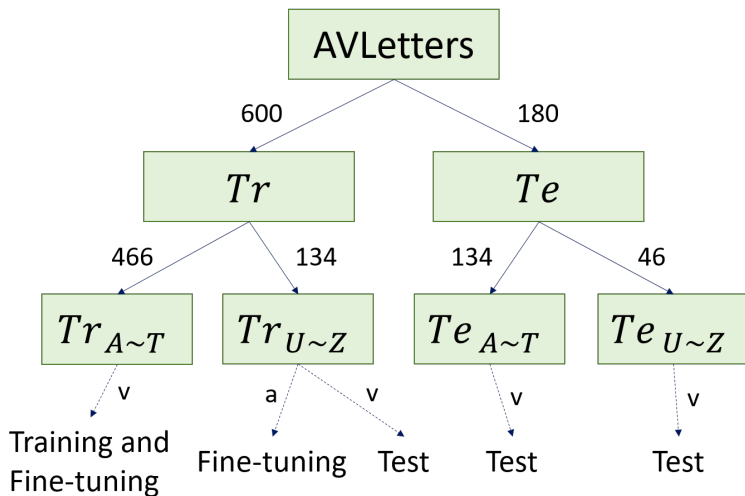


Results – Shared Representation

- Shared: Perceptron that exploits the shared representation learned by a bimodal CAE.

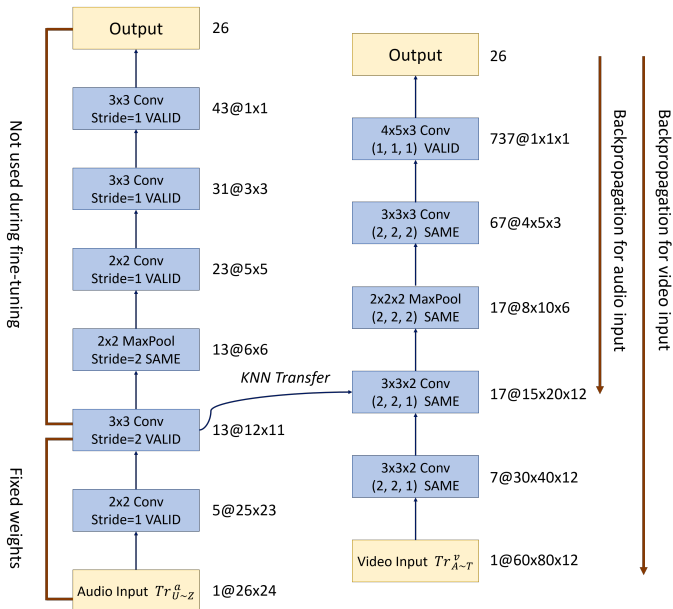
		Raw	CAE features	CAE architecture	Shared
Intensity	train	69.47 %	78.87 %	91.29 %	85.85 %
	test	32.64 %	50.24 %	65.44 %	53.38 %
Depth	train	63.64 %	79.61 %	88.80 %	81.83 %
	test	29.93 %	41.64 %	55.62 %	42.85 %

AVSR Knowledge Transfer



Pretrained Audio Network

Pretrained Video Network



AVSR Knowledge Transfer

Fine-tuned for 160 steps.

For **Exp1**, **Exp2** and **Exp3**, we have respectively $\alpha_0 = 0.001, 0.005, 0.001$ and $p_a = 0.85, 0.85, 1$.

Notice that since $p_a = 1$ for **Exp3** no video data is given in input during fine-tuning.

	Tr^v	$Tr_{A \sim T}^v$	$Tr_{U \sim Z}^v$	Te^v	$Te_{A \sim T}^v$	$Te_{U \sim Z}^v$
No transfer	77.67 %	100 %	0 %	40.56 %	54.48 %	0 %
Exp1	81.17 %	98.28 %	21.64 %	39.44 %	47.76 %	15.22 %
Exp2	40.83 %	51.07 %	5.22 %	23.89 %	30.60 %	4.35 %
Exp3	19.67 %	12.23 %	45.52 %	12.22 %	2.24 %	41.34 %

Conclusion

- Datasets, hyperparameters
- Applications in robotics
- Other approaches