



# Chapter 7, Part 1: Introduction to Deep Learning

Advanced Topics in Statistical Machine Learning

---

Tom Rainforth

Hilary 2024

[rainforth@stats.ox.ac.uk](mailto:rainforth@stats.ox.ac.uk)

<https://youtu.be/cQ54GDm1eL0>

# Motivation: High Dimensional Inputs

Many common data types are (very) high dimensional, e.g. text, images, and videos



**Disentangling Disentanglement in Variational Autoencoders**

---

Emile Mathieu<sup>\*1</sup>, Tom Rainforth<sup>\*1</sup>, N. Siddharth<sup>\*2</sup>, Yee Whye Teh<sup>1</sup>

---

**Abstract**

We develop a generalisation of disentanglement in variational autoencoders (VAEs)—decomposition of the latent representation—characterising it as the fulfilment of two factors: a) the latent encodings of data have an appropriate level of overlap, and b) the aggregate encoding of the data conforming to a desired structure, represented through the decomposition permits disentanglement, i.e., capturing independent latent factors, as a special case, but also allows for a much richer class of properties to be imposed on the latent representation, such as sparsity, clustering, and submodularity, which capture more complex hierarchical dependency relationships. We show that the  $\beta$ -VAE varies from the standard VAE predominantly in its control of latent overlap and that the proposed decomposition approach, at no additional prior, is objective in invariant the invariance with respect to the choice of latent factor. We demonstrate disentanglement with little or no detriment to reconstructions. We further demonstrate how other choices of prior can result in producing different decompositions and introduce an alternative training objective that allows the control of both decomposition factors in a principled manner.

**1. Introduction**

An oft-stated motivation for learning disentangled representations of data with deep generative models is a desire to achieve ‘invariance’ (Higgins et al., 2017; Bell et al., 2017)—particularly the ‘incomposability’ (see §3.2 in Lipton, 2016) of latent representations to allow a disentangling objective that allows the control of both decomposition factors in a principled manner.

**2. Related work**

Independence factors of variation (Almeida et al., 2017; Assari and Soh, 2019; Burgess et al., 2018; Chen et al., 2018, 2017; Eastwood and Williams, 2018; Esmaili et al., 2019; Higgins et al., 2017; Kar and Muth, 2018; Higgins et al., 2018; Zhao et al., 2017), typically evaluating this using purpose-built synthetic data (Eastwood and Williams, 2018; Higgins et al., 2018; Kar and Muth, 2018), where generative factors are independent of one another.

This conventional view of disentanglement, as recovering independence, has subsequently motivated the development of formal evaluation metrics for independence (Eastwood and Williams, 2018; Kar and Muth, 2018), which in turn have driven the development of methods to optimise these metrics, often by employing regularisers explicitly encouraging independence in the representations (Eastwood and Williams, 2018; Esmaili et al., 2019; Kar and Muth, 2018). We argue that such an approach is non-generalisable, and potentially counterproductive, to learning interesting representations that are more complex than priors where such independence representations cannot accurately mimic the generation of high dimensional data from low dimensional latent spaces, and more richly model dependencies are required.

We propose a generalisation of disentanglement in VAEs—disentangling their latent representations—which helps avoid such pitfalls. We characterise disentanglement in VAEs as the fulfilment of two factors: a) the latent encodings of data having an appropriate level of overlap, and b) the aggregate encoding of data following a desired structure represented through the prior. We emphasise that neither of these factors is sufficient in isolation: without an appropriate level of overlap, encodings can degenerate to a lookup table, and without an appropriate level of structure, the encodings will not be disentangled, and without the aggregate encoding of data following a desired structure, the encodings do not decompose as desired.

Disentanglement implicitly makes a choice of decomposition that the latent features are independent of one another. We introduce *explicit* and *envelope* priors to provide incentives to learn disentangled representations through judicious choices of structure in the prior, and to introduce a more general framework flexible enough to capture alternate, more complex, forms of disentanglement such as sparsity, clustering, hierarchical structuring, or independent subspaces.

---

<sup>\*</sup>Equal contribution. Department of Statistics, Department of Engineering, University of Oxford. Correspondence to: Emile Mathieu (emile.mathieu@maths.ox.ac.uk), Tom Rainforth (rainforth@stats.ox.ac.uk), N. Siddharth (n.siddharth@robots.ox.ac.uk).

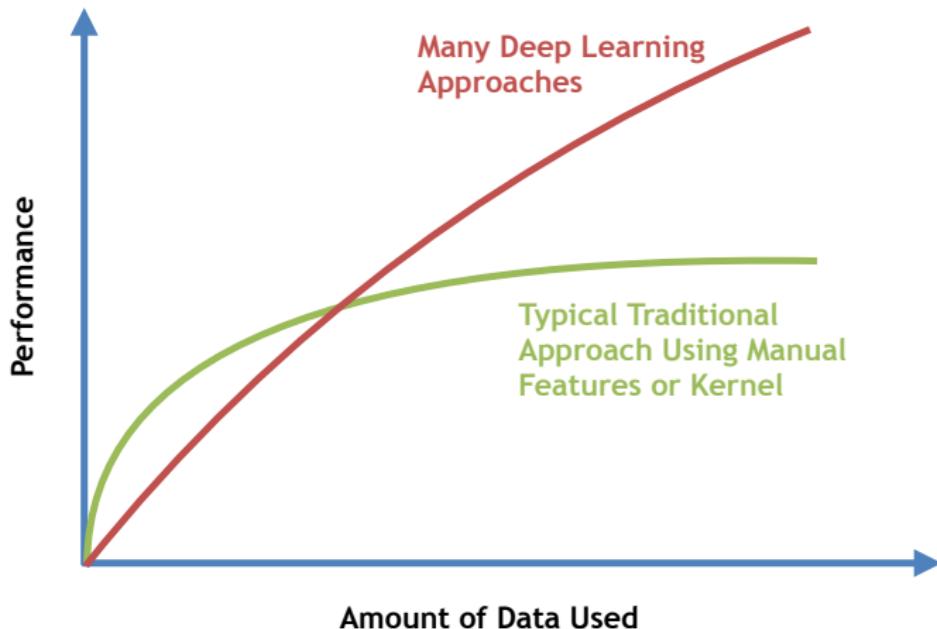
Proceedings of the 37<sup>th</sup> International Conference on Machine Learning, Long Beach, California, PMLR 97, 2019. Copyright 2019 by the authors.

## Motivation: Flexibility and Complicated Output Spaces



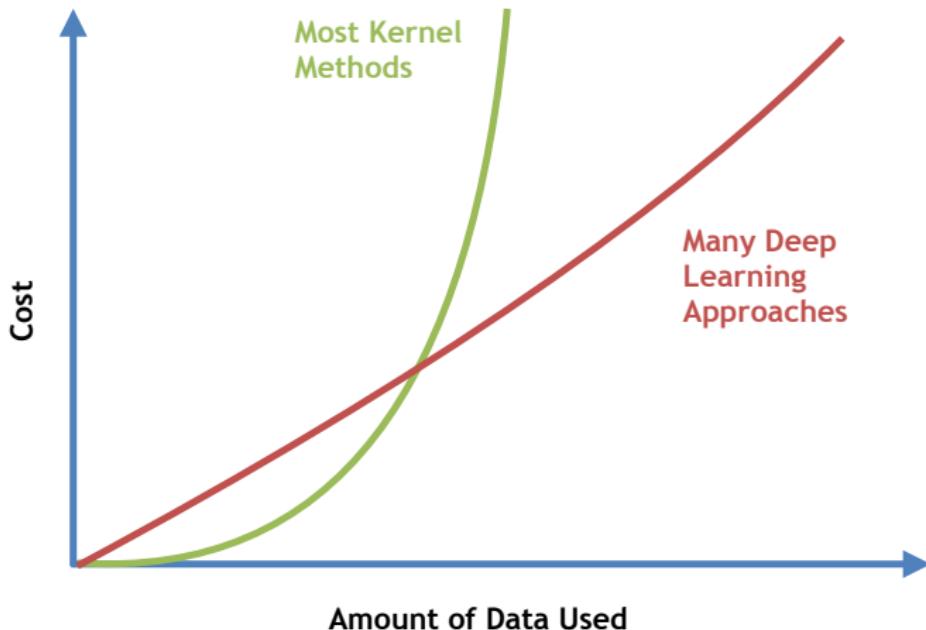
**Figure 1:** Fake Rembrandt generated by a deep learning algorithm.  
Image credit: <https://www.nextrembrandt.com/>

## Motivation: Performing Well on Huge Datasets



**Figure 2:** Characterization of deep learning method's performance as training dataset size increases (somewhat overgeneralized)

# Motivation: Graceful Computational Scaling



**Figure 3:** Characterization of deep learning method's cost as training dataset size increases (somewhat overgeneralized)

# Feature Engineering

Even though they imply an infinite number of features, kernel methods still require an effective similarity function between inputs and that requires specialist knowledge about the space of inputs

- To know if inputs are similar, we need to know what about them is important

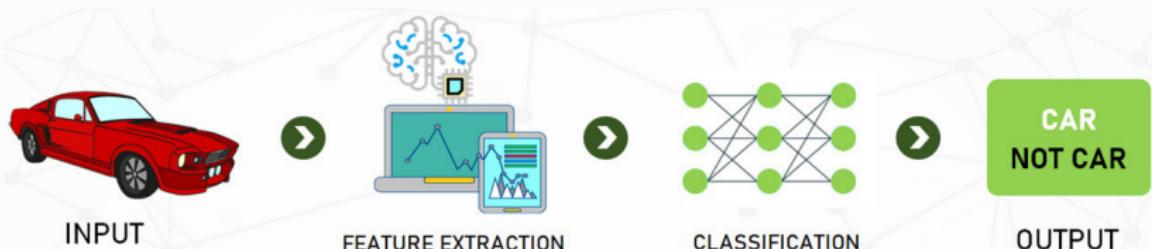
We can think of this problem as that of **feature engineering**: to perform effective machine learning we need effective **representations** of our inputs that encapsulate salient information

- We need to pick out **patterns** in the space of inputs that allow us to assess similarities and differences between datapoints
- We want to **ignore** unimportant aspects and noise

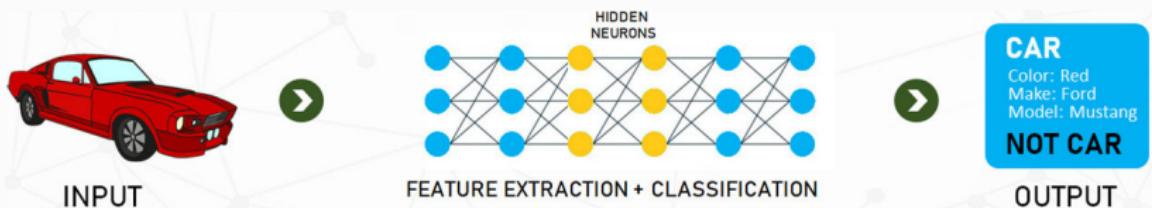
**Manually** constructing effective features (or kernels) is extremely difficult in **high dimensions**

# End-To-End Learning

## Classical Feature Engineering Approach



## End-To-End Deep Learning Approach



**Figure 4:** Deep learning looks to replace manual feature engineering with an end-to-end approach that learns everything at the same time

## End-To-End Learning (2)

A classical feature-engineering approach will take a fixed feature mapping  $\varphi : \mathcal{X} \rightarrow \mathbb{R}^m$  and then perform ERM to optimize the parameters of a predictive model  $f_\theta : \mathbb{R}^m \rightarrow \mathcal{Y}$

$$\theta^* = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n L(y_i, f_\theta(\varphi(x_i))) + r(\theta)$$

where  $L$  is our loss function and  $r$  is a regularizer.

The starting point for deep learning is to instead take an end-to-end approach where we also parameterize the feature mapping  $\varphi_\eta : \mathcal{X} \rightarrow \mathbb{R}^m$  and then optimize all parameters

$$\{\theta^*, \eta^*\} = \arg \min_{\theta, \eta} \frac{1}{n} \sum_{i=1}^n L(y_i, f_\theta(\varphi_\eta(x_i))) + r(\theta, \eta)$$

We are thus replacing manual feature engineering with **feature learning** that is performed simultaneously to training the model

## Using Multiple Layers

The key idea of deep learning is now to take this a step further and introduce hierarchical **composition** of **nonlinear** functions

We do this by replacing  $f_\theta(\varphi_\eta(x))$  with

$$f_{\theta^m}^m \left( f_{\theta^{m-1}}^{m-1} \left( \dots f_{\theta^2}^2 \left( f_{\theta^1}^1(x) \right) \dots \right) \right)$$

such that our overall predictive model is  $f_\theta : \mathcal{X} \mapsto \mathcal{Y}$  is given by

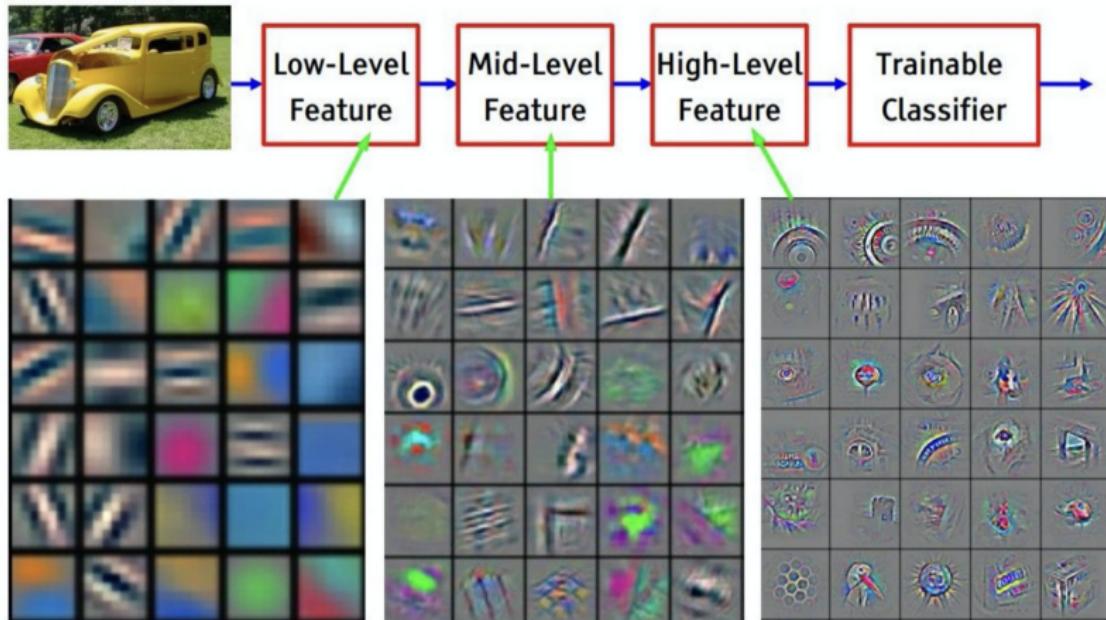
$$f_\theta = f_{\theta^m}^m \circ f_{\theta^{m-1}}^{m-1} \circ \dots \circ f_{\theta^2}^2 \circ f_{\theta^1}^1$$

We can also think of this as having multiple **layers** of features,  $h^\ell \in \mathbb{R}^{d_\ell}$ , such that  $h^0 = x$ ,  $h^m = f_\theta(x)$  and

$$h^\ell = f_{\theta^\ell}^\ell \left( h^{\ell-1} \right)$$

“**Deep**” now refers to having many such layers, where we **learn** the parameters  $\theta^\ell$  for each of the associated mappings

# Deeper Layers Represent More Complex Features



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

**Figure 5:** Visualization of features represented by different layers.

Individual images in the grid correspond to example dimensions of a  $h^\ell$ .

## Desiderata for $f_\theta$

To more precisely understand why taking a deep approach is beneficial, let's first consider the desirable properties for the class of predictors  $f_\theta$ :

- **Powerful:** we want  $f_\theta$  to be able to accurately capture complex input–output relationships
- **Flexible and Expressive:** we want to be able to apply the approach to a wide array of problems and use customized setups that are tailored to specific tasks
- **Easily Trainable and Scalable:** we need to be able to effectively and tractably train good predictors, even with large amounts of high dimensional training data.

## Achieving these Desiderata

- To be sufficiently powerful with large quantities of data, we need to be able to use a **large number of parameters**  $\theta$
- To allow effective training in high dimensions we will typically need to be able to take **derivatives**  $\nabla_{\theta} f_{\theta}(x)$
- To be flexible and expressive, we will generally require complex hypothesis classes to be constructible by piecing together simpler building blocks

**Deep learning** is a compositional framework that achieves these by building up complex hypothesis classes from compositions of simple **differentiable** building blocks

A particular configuration of these building blocks represents a hypothesis class and is known as an **architecture**

# Deep Learning

Deep learning is a machine learning approach where we construct a parameterized predictive function using a computational graph of simpler differentiable functional blocks, and then train this on data using some form of gradient-based optimization.

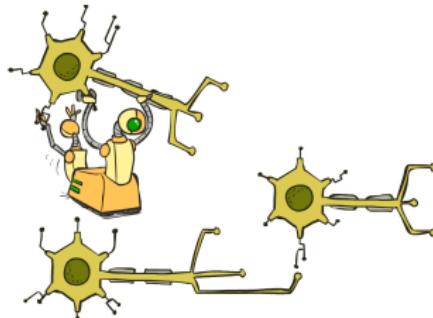


Image credit: Dan Klein and Pieter Abbeel

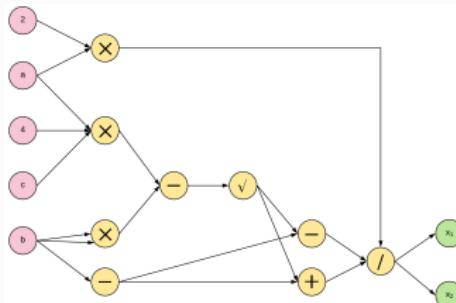


Image credit: Ekaba Bisong

# Differentiable Programming



**Yann LeCun**

January 5 at 10:13pm ·

...

OK, Deep Learning has outlived its usefulness as a buzz-phrase.  
Deep Learning est mort. Vive Differentiable Programming!

Yeah, Differentiable Programming is little more than a rebranding of the modern collection Deep Learning techniques, the same way Deep Learning was a rebranding of the modern incarnations of neural nets with more than two layers.

But the important point is that people are now building a new kind of software by assembling networks of parameterized functional blocks and by training them from examples using some form of gradient-based optimization.

## Compositions of Differentiable Functions are Differentiable

The reason we can do this is because the chain rule ensures that the composition of differentiable functions is itself differentiable:

$$\frac{\partial L(y_i, f_\theta(x_i))}{\partial \theta^k} = \frac{\partial L(y_i, f_\theta(x_i))}{\partial h^m} \frac{\partial h^m}{\partial h^{m-1}} \cdots \frac{\partial h^{k+1}}{\partial h^k} \frac{\partial h^k}{\partial \theta^k} \Big|_{x=x_i}$$

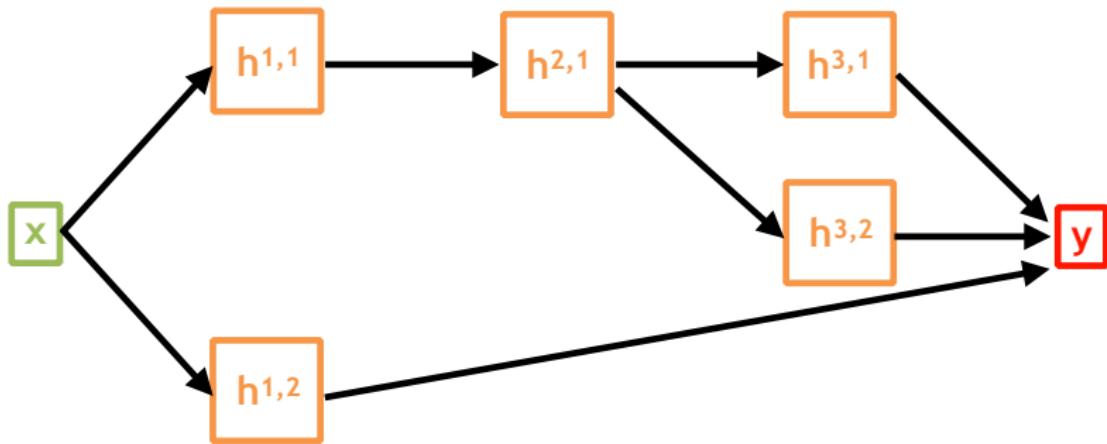
where each  $\frac{\partial h^\ell}{\partial h^{\ell-1}}$  is a  $d_\ell \times d_{\ell-1}$  matrix representing the Jacobian of  $f_{\theta^\ell}$  with respect to its inputs.

Note that, except for the first layer, we require that the  $f_{\theta^\ell}$  are differentiable with respect to their parameters **and** their inputs

This also holds for general differentiable **computation graphs** by allowing individual  $h^\ell$  to be split up into separate nodes

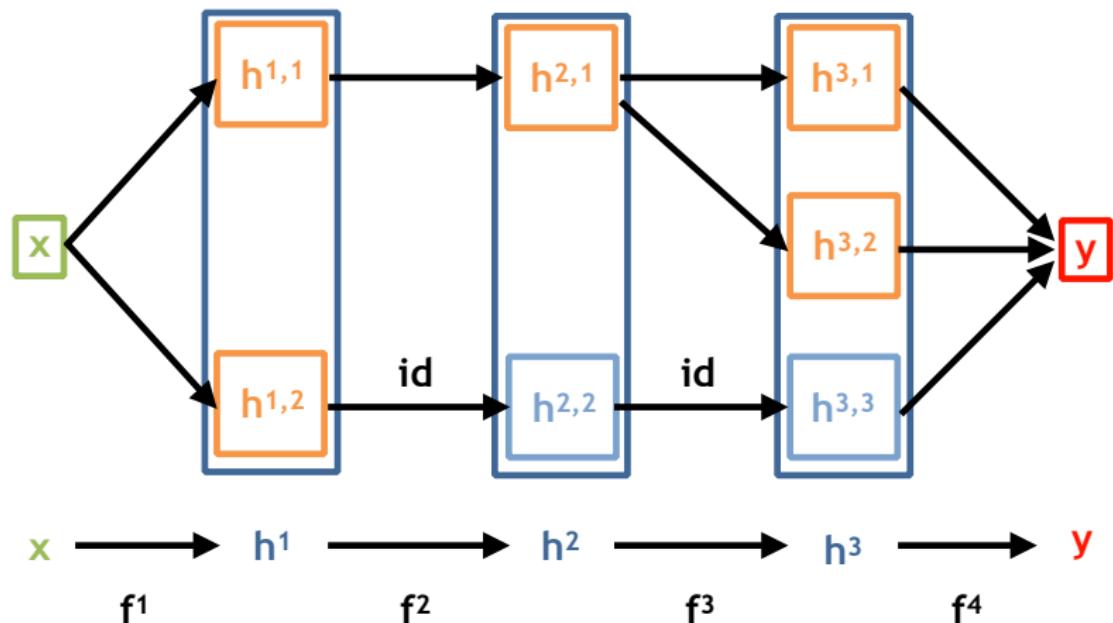
$h^\ell = [h^{\ell,1}, h^{\ell,2}, \dots, h^{\ell,t_\ell}]$ , or equivalently grouping nodes together as layers and padding with identity mappings as required

# Computation Graphs



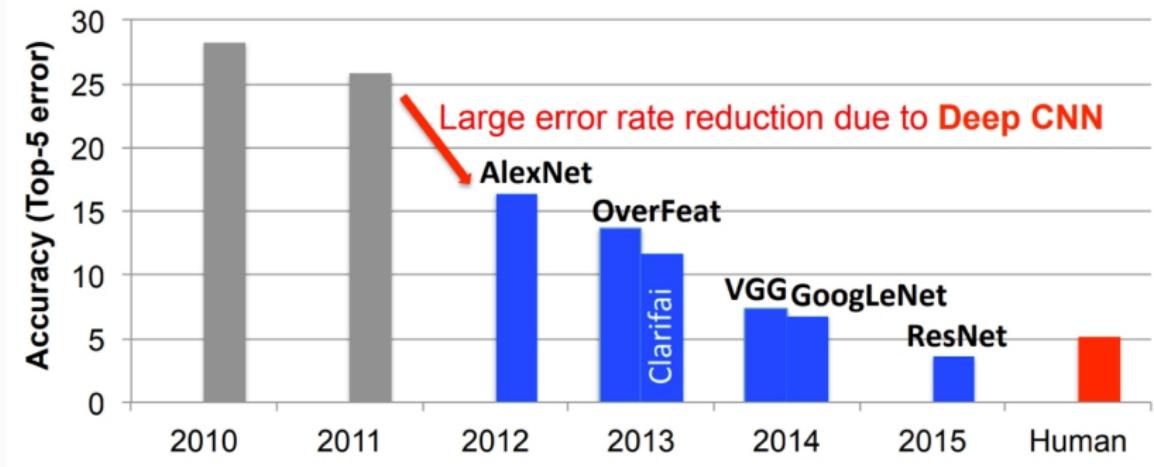
**Figure 6:** Example computation graph and how this can be represented as series of layers by grouping and padding nodes.  $id$  represents the identity function  $f(x) = x$ .

# Computation Graphs



**Figure 6:** Example computation graph and how this can be represented as series of layers by grouping and padding nodes. id represents the identity function  $f(x) = x$ .

## Applications: Image Classification



**Figure 7:** State of the art performance on ImageNet (a large image classification dataset) over time. Figure taken from Russakovsky et al IJCV 2015.

# Applications: Image Captioning



"man in black shirt is playing guitar."



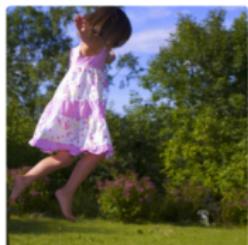
"construction worker in orange safety vest is working on road."



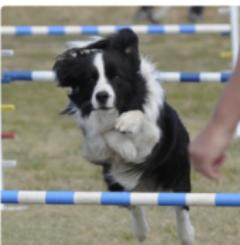
"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"girl in pink dress is jumping in air."



"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."



"man in blue wetsuit is surfing on wave."

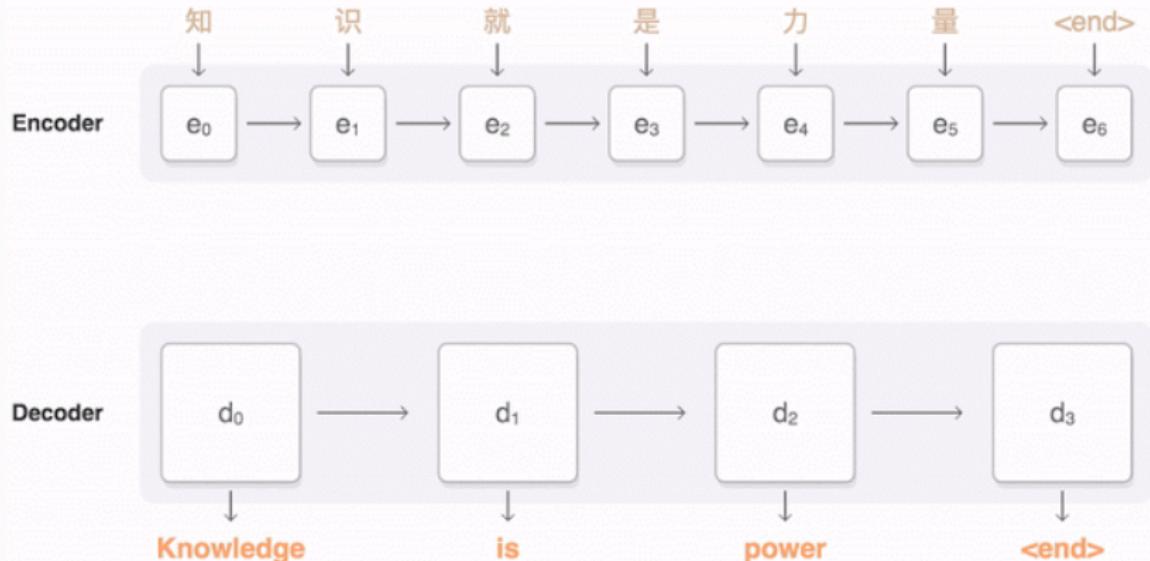
Karpathy & Fei-Fei, 2015; Donahue et al., 2015; Xu et al, 2015; many more

**Figure 8:** MS COCO Image captioning challenge. Image credit Dan Klein and Pieter Abbeel

## Applications: Deep Fakes

[https://youtu.be/h\\_jrebvmPlk](https://youtu.be/h_jrebvmPlk)

## Applications: Test Translation



**Figure 9:** Translating Cantonese into English. Image Credit: Google AI Research Blog

## Applications: Large Language Models

<https://chat.openai.com/chat>

# Applications: Learning to Play Games

Learning to play Mario Kart

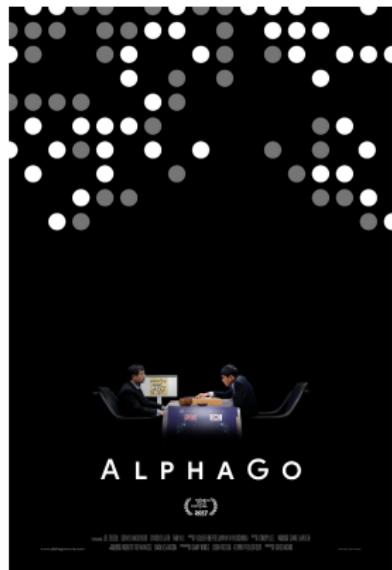
[https://www.youtube.com/watch?v=Tnu40\\_xEmVk](https://www.youtube.com/watch?v=Tnu40_xEmVk)

Beating top players at Starcraft

<https://youtu.be/DMXvkbAtHNY>

Learning to park (shows requirement for a lot of training/simulated data)

[https://youtu.be/VMp6pq6\\_QjI](https://youtu.be/VMp6pq6_QjI)



**Figure 10:** Beating the world champion Go player