

# Deep Learning & Applications

**Michalis Vazirgiannis**

**Data Science & Mining group**  
LIX, Ecole Polytechnique,

DASCIM web page: <http://www.lix.polytechnique.fr/dascim>

Google Scholar: <https://bit.ly/2rwmvQU>

Twitter: @mvazirg

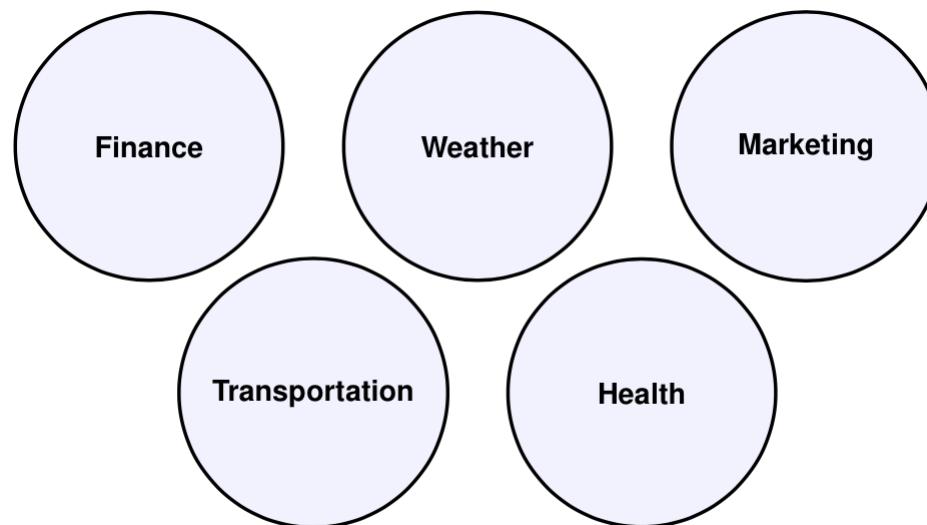
**November 2022**

- Time Series
- Recommendations
- Biomed

# TIME SERIES - Motivation

*Why focus on Times Series models?*

- ① Constant need to predict the future.
- ② Focus on minimizing risk, developing long-term strategies, making decisions instantly.
- ③ Availability of a great amount of data.
- ④ Impact on a huge range of applications.



# A Time Series Example - Stock prices

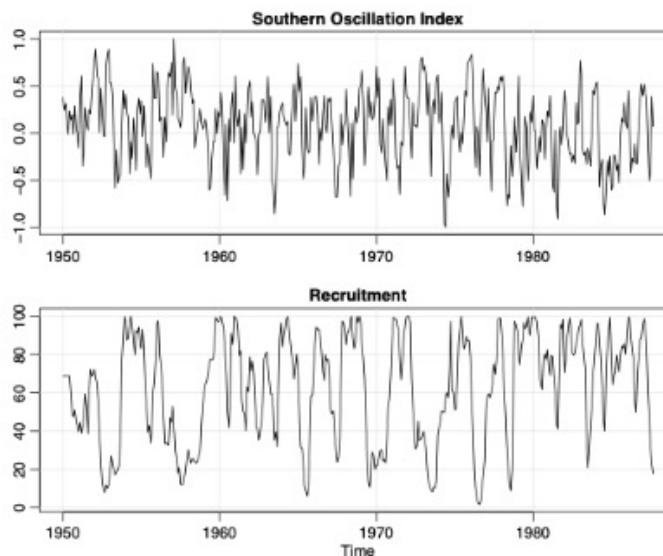
APPLE, TESLA stocks, A 5-year weekly evolution<sup>1</sup>



<sup>1</sup>source: [finance.yahoo.com](https://finance.yahoo.com)

# A Time Series Example - weather cycles and fish

Monthly SOI and Recruitment (estimated new fish), 1950-1987 <sup>4</sup>

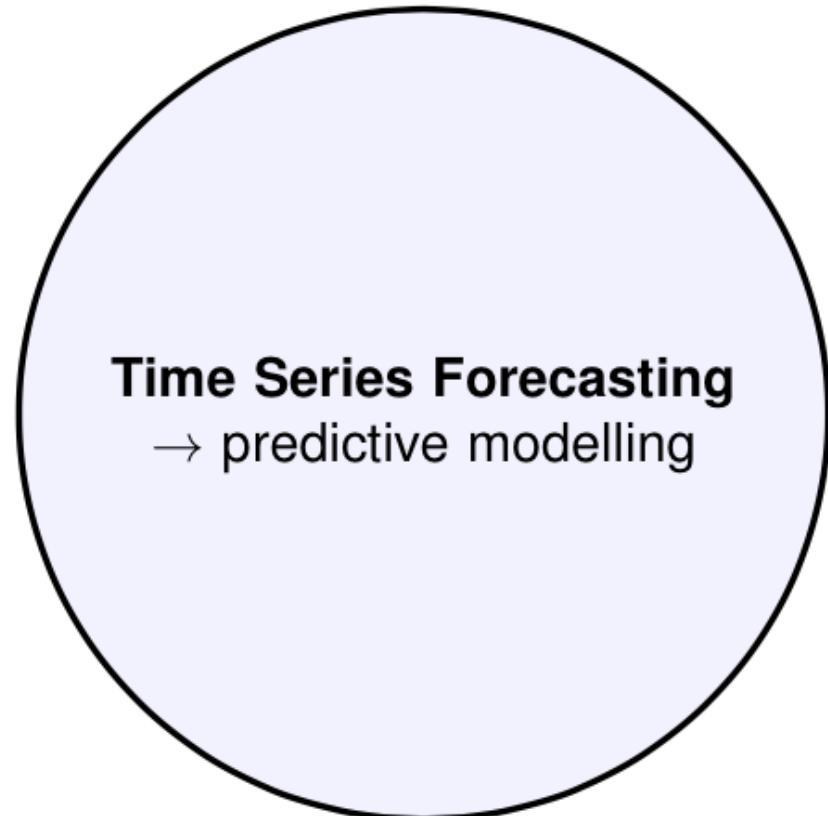
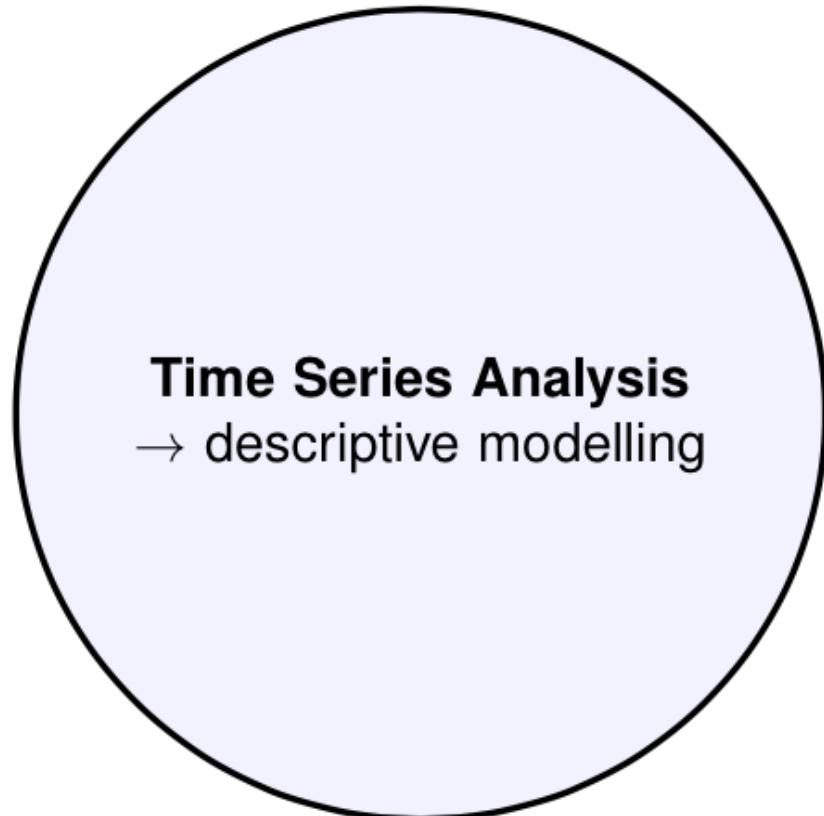


- SOI: changes in air pressure, related to sea surface temperatures - central Pacific Ocean.
- Central Pacific warms every 3-7 years (El Niño effect blamed for global extreme weather).
- Both series exhibit repetitive behavior, (cycles) - underlying processes: rate or frequency of oscillation

---

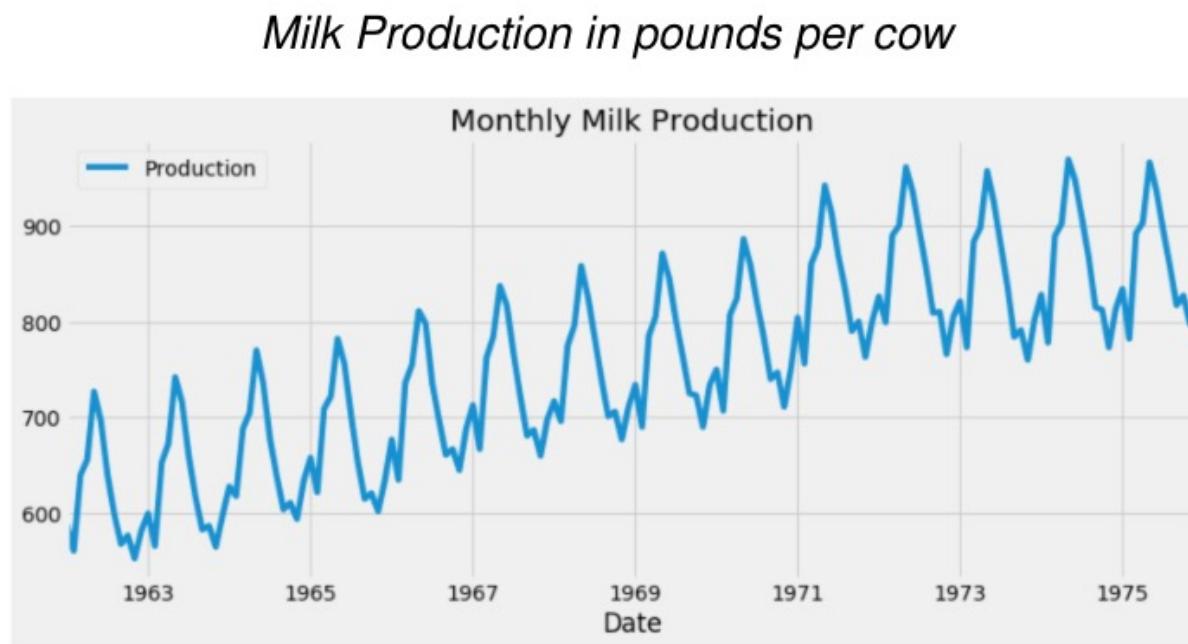
<sup>4</sup> R. H. Shumway, D. S. Stoffer, *Time Series Analysis and Its Applications, With R Examples*. Springer Texts in Statistics

# Time Series - Domains of study



# Components of a Time Series

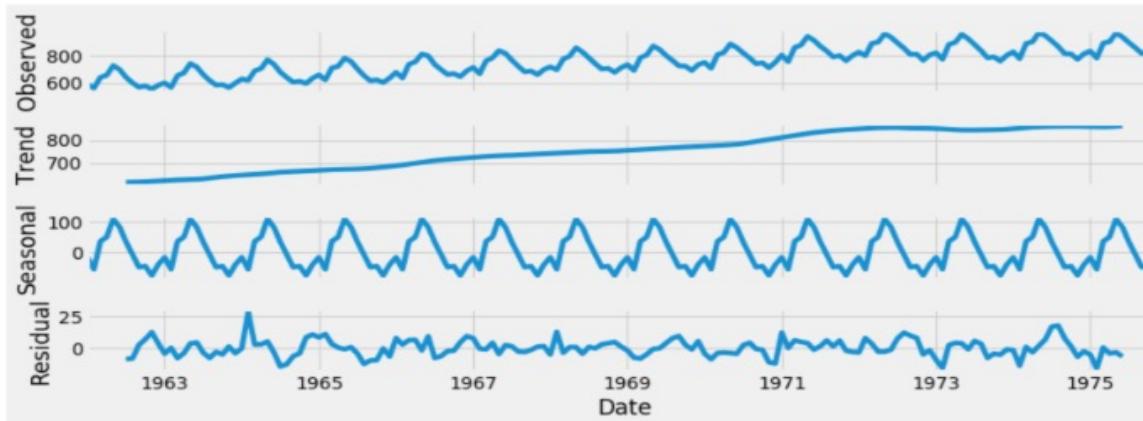
The first step in analyzing a time series for a predictive model is to identify the underlying pattern of the data over time.



In general, a time series is affected by four components, i.e. **trend, seasonal, cyclical and irregular components**.

# Components of a Time Series

*Time Series Decomposition - Monthly Milk Production*



- **Trend :**

The long-term gradual change in the series. This is the simplest trend pattern, as it demonstrates long-term growth or decline.

- **Seasonality:**

Predictable, short-term patterns that occur within a single unit of time and repeat indefinitely.

- **Noise or irregular variation (residual)**

Random variation not regular and also do not repeat in a particular pattern, caused by incidences such as war, strike, earthquake etc.

# Time Series Decomposition

Based on these four components, a time series can be decomposed using the following approaches. Decomposition is primarily used for time series analysis, and as an analysis tool it can be used to inform forecasting models on your problem.

## ① Additive Model:

- $Y(t) = T(t) + S(t) + I(t)$
- These three components are independent of each other.
- Assumption: seasonal variation is constant over time.

## ② Multiplicative Model:

- $Y(t) = T(t) \times S(t) \times I(t)$
- These three components are not necessarily independent and they can affect one another.
- Assumption: seasonal variation *increases* over time.

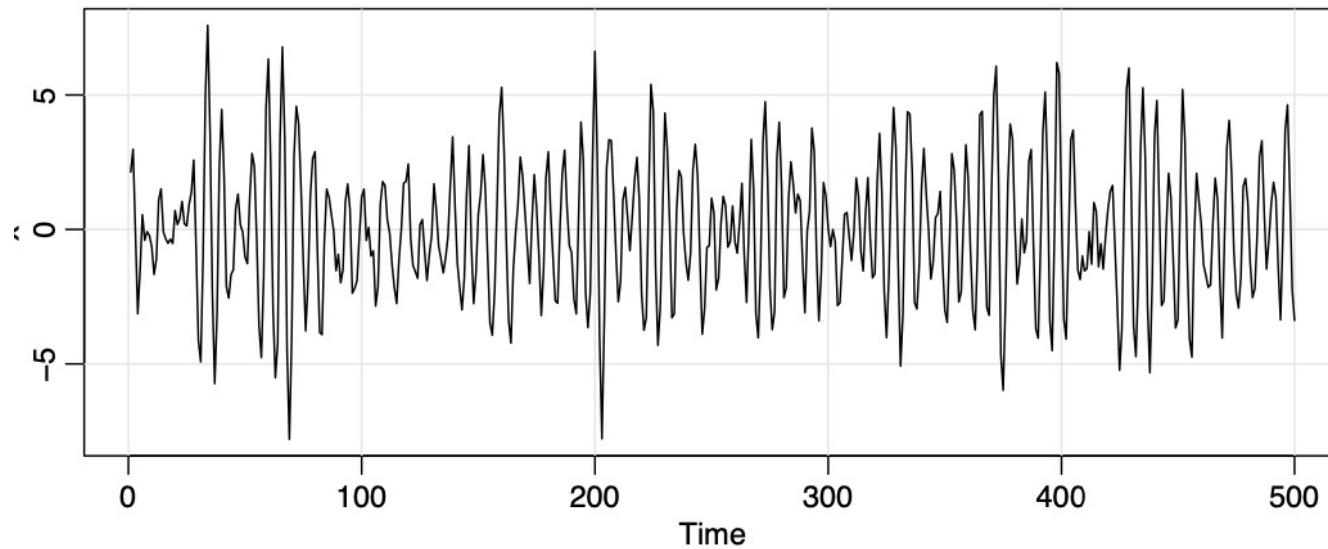
where  $T(t)$ : trend (with cycles),  $S(t)$ : seasonal,  $I(t)$ : irregular (random) components

# Time series as a stochastic process

- We can define a time series as a collection of random variables indexed based on the order they are obtained in time,  $X_1, X_2, X_3, \dots, X_t$  will typically be discrete and vary over the integers  $t = 0, \pm 1, \pm 2, \dots$
- The collection of the random variables  $X_t$  is referred to as a stochastic process, while the observed values are referred to as a realization of the stochastic process.
- A time series observed as **a collection of  $n$  random variables** at arbitrary time points  $t_1, t_2, \dots, t_n$ , for any positive integer  $n$ , is provided by the joint distribution function, using the  $n$  constants,  $c_1, c_2, \dots, c_n$  i.e.:
  - $F_{t_1, t_2, \dots, t_n}(c_1, c_2, \dots, c_n) = Pr(X_{t_1} \leq c_1, X_{t_2} \leq c_2, \dots, X_{t_n} \leq c_n)$

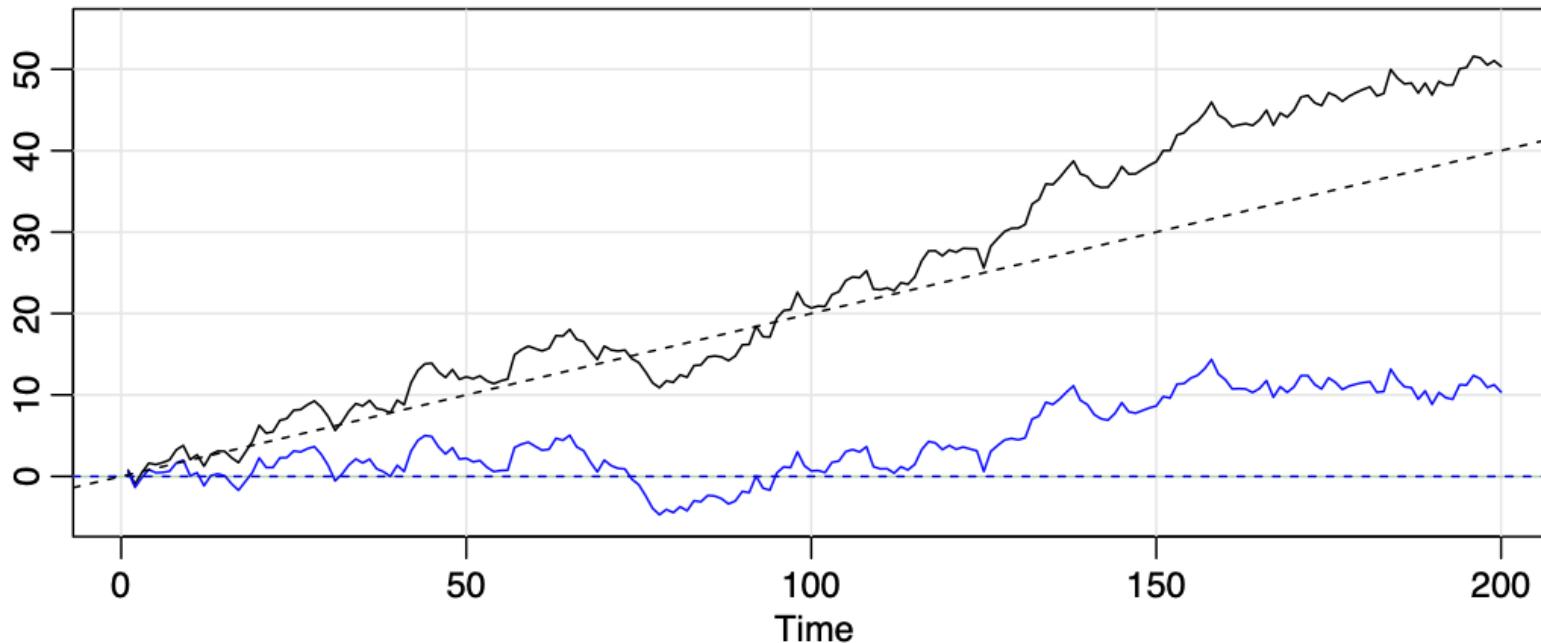
# Autoregressive time series

- Assume white noise series  $w_t$  :  $x_t = x_{t-1} - .9x_{t-2} + w_t$



# Random Walk with Drift time series

- *Time series model:*  $x_t = \delta + x_{t-1} + w_t$   $t = 1, 2, \dots, x_0 = 0$ , and  $w_t$  is white noise.  $\delta$ : constant is called the *drift*
- when  $\delta = 0$ , is called simply a *random walk*.

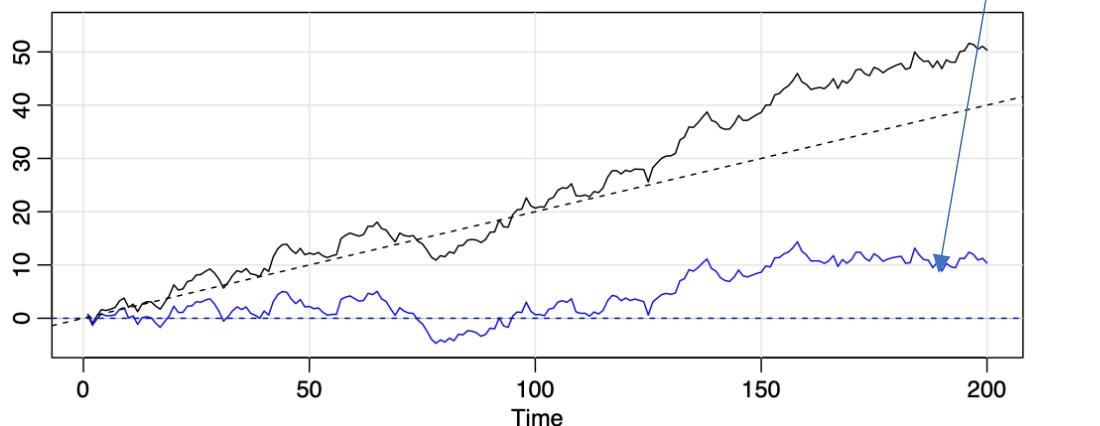


# Mean Function - Time Series

- The mean function is defined as  $\mu_t = \mu_{X_t} = E[X_t] = - \int_{-\infty}^{\infty} xf_t(x)dx$ , provided it exists, where  $E$  the usual expected value operator.
- For White Noise:  $\mu_{w_t} = E(w_t) = 0$
- For Random Walk:

$$\mu_{X_t} = E[X_t] = E[X_{t-1} + \delta + w_t] = E[\delta_t + w_t] = \delta_t + \sum_{i=1}^t E[w_i]$$

$$\Rightarrow \mu_{X_t} = \delta_t$$



# Autocovariance - Time Series

- **Autocovariance Function**

- is defined as the second moment product

$$\gamma(s, t) = \gamma_X(s, t) = \text{cov}(X_s, X_t) = E[(X_s - \mu_s)(X_t - \mu_t)], \forall s, t$$

- $\gamma(s, t) = \gamma(t, s), \forall s, t$

- The autocovariance **measures the linear dependence between two points** on the same series observed at different times.

- **cross-covariance function** *between two series,  $x_t$  and  $y_t$ , is*

$$\gamma_{xy}(s, t) = \text{cov}(x_s, y_t) = E[(x_s - \mu_{xs})(y_t - \mu_{yt})].$$

# Auto correlation - Time Series

- **Autocorrelation Function (ACF)**

- $\rho(s, t) = \frac{\gamma(s, t)}{\sqrt{\gamma(s, s)\gamma(t, t)}}$

- From Cauchy-Schwarz inequality

$$|\gamma(s, t)|^2 \leq \gamma(s, s)\gamma(t, t) \Rightarrow -1 \leq \rho(s, t) \leq 1$$

- ACF measures the **linear predictability of  $X_t$  using only  $X_s$**

- If we can predict  $X_t$  perfectly from  $X_s$  through a linear relationship, then ACF will be either +1 or -1.

- *The cross-correlation function (CCF) between two series,  $x_t$  and  $y_t$*

$$\rho_{xy}(s, t) = \frac{\gamma_{xy}(s, t)}{\sqrt{\gamma_x(s, s)\gamma_y(t, t)}}.$$

# Stationary Time Series

- **strictly stationary time series:** the probabilistic behavior of every collection of values  $\{x_{t_1}, x_{t_2}, \dots, x_{t_k}\}$  is identical to that of the time shifted set  $\{x_{t_1+h}, x_{t_2+h}, \dots, x_{t_k+h}\}$  for all  $h$  or:

$$\Pr\{x_{t_1} \leq c_1, \dots, x_{t_k} \leq c_k\} = \Pr\{x_{t_1+h} \leq c_1, \dots, x_{t_k+h} \leq c_k\}$$

- **weakly stationary time series**
  - The time series  $X_t, t \in \mathbb{Z}$  is said to be weak stationary if:
    - ➊  $E[x_t^2] < \infty, \forall t \in \mathbb{Z}$
    - ➋  $E[x_t] = \mu, \forall t \in \mathbb{Z}$
    - ➌  $\gamma_X(s, t) = \gamma_X(s + h, t + h), \forall s, t, h \in \mathbb{Z}$
  - *mean value function,  $\mu_t$ , is constant and does not depend on time  $t$*
  - *autocovariance function,  $\gamma(s, t)$  depends on  $s$  and  $t$  only through their difference  $|s - t|$ .*
  - Stationary datasets are those that have a *stable mean and variance*,

# AR Models

**Autoregressive Models (AR Models)** express the current value of the series,  $X_t$  as a linear combination of  $p$  past values,  $X_{t-1}, X_{t-2}, \dots, X_{t-p}$ , including a random error in the approximation.

**Definition:** An autoregressive model of order  $p$ , **AR(p)**, is defined as

$$X_t = \phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} + w_t \Rightarrow$$

$$X_t = \sum_{i=1}^p \phi_i X_{t-i} + w_t$$

where  $X_t$  is stationary,  $\phi_1, \phi_2, \dots, \phi_p$  are model parameters (where  $p$  the length of historical points) and  $w_t \sim wn(0, \sigma_w^2)$ .

# AR Models: AR(0), AR(1), AR(p)

- ①  $AR(0)$  is the simplest AR process and has no dependence between the terms (white noise).
- ②  $AR(1)$  is defined as  $X_t = \phi_1 X_{t-1} + w_t$ 
  - if  $|\phi_1| \approx 0$  the process behaves like white noise.
  - if  $\phi_1 = 1$  the process is equivalent to random walk with infinite variance (dependent on  $t$ , non-stationary).
  - if  $\phi_1 < 0$  the process swings between positive and negative values.
- ③ For the general  $AR(p)$  process:
  - if  $h > p$  the PACF is
$$\phi_{hh} = \text{corr}(X_{t+h} - \hat{X}_{t+h}, X_t - \hat{X}_t) = \text{corr}(w_{t+h}, X_t - \hat{X}_t) = 0.$$
  - if  $h \leq p$  then  $\phi_{pp}$  is not 0.
  - thus the identification of an AR model can be efficiently done with the PACF.

# AR Models: Parameters Estimation

- ①  $p$  is a hyperparameter for the  $AR(p)$  process, thus when fitting an  $AR(p)$  model  $p$  is known and we focus on estimating the coefficients  $(\phi_1, \phi_2, \dots, \phi_p)$ .
- ② Coefficients estimation can be done by different approaches:
  - Maximum Likelihood Estimation estimator (MLE).
  - Ordinary Least Squares estimator (OLS).

# Moving Average models

One problem of AR model is the ignorance of correlated noise structures in the time series. In other words,  $w_t$  and its historical terms  $w_{t-1}, w_{t-2}, \dots, w_{t-q}$ , include information that can be utilized for predictive models.

**Definition:** A Moving Average Model (MA) of order  $q$ , **MA( $q$ )**, is defined as

$$X_t = \theta_1 w_{t-1} + \theta_2 w_{t-2} + \dots + \theta_q w_{t-q} + w_t \Rightarrow$$

$$X_t = \sum_{j=1}^q \theta_j w_{t-j} + w_t$$

where  $X_t$  is stationary,  $\theta_1, \theta_2, \dots, \theta_p$  are model parameters (where  $q$  the length of historical points) and  $w_t \sim \text{wn}(0, \sigma_w^2)$ .

- looks like a regression model, the difference is that the  $w_t$  is not observable.
- In contrast with AR models, finite MA models are always stationary (observation is a weighted moving average over past forecast errors).

# ARMA Models

Autoregressive and Moving Average models can be combined and form **ARMA models**.

**Definition:** A stationary time series is **ARMA(p, q)** if

$$X_t = w_t + \sum_{i=1}^p \phi_i X_{t-i} + \sum_{j=1}^q \theta_j w_{t-j}$$

where  $\phi_p, \theta_q \neq 0$  and  $w_t \sim wn(0, \sigma_w^2)$ .

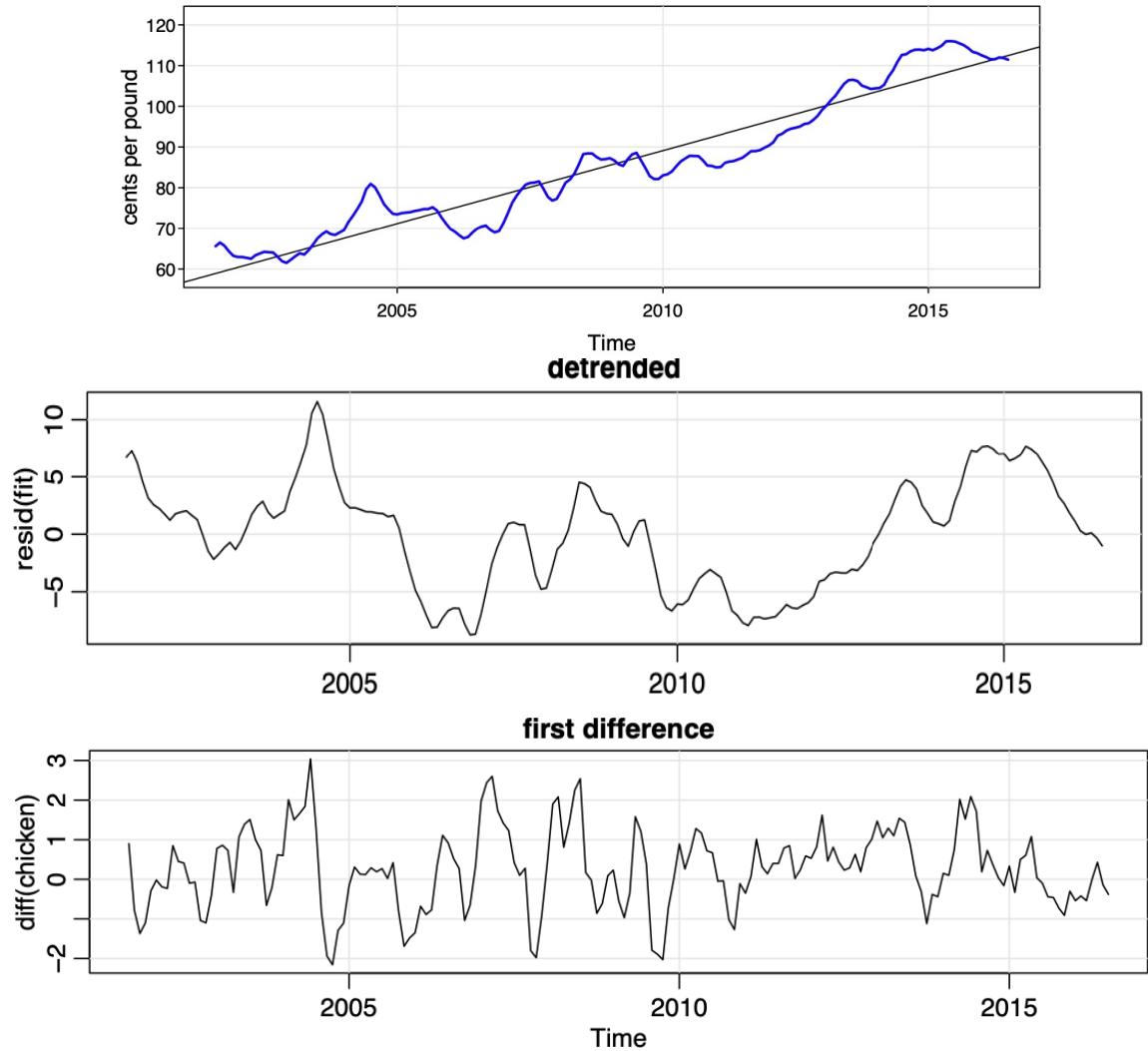
# Transforming Time Series to Stationary

One prerequisite of ARMA models is the stationarity of time series

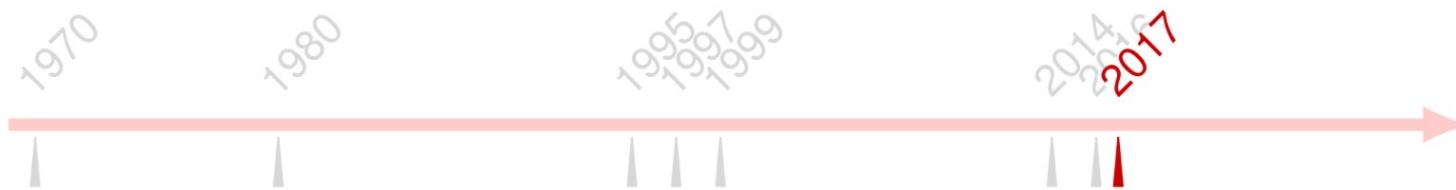
- Assume non-stationary time series:  $x_t = \mu_t + y_t$
- $x_t$  are the observations,  $\mu_t$  trend,  $y_t$  is a stationary process.
- Remove the trend: estimate of the trend component:  $\hat{\mu}_t$ , and then work with the residuals

$$\hat{y}_t = x_t - \hat{\mu}_t$$

- Differencing:  $x_t = x_t - x_{t-1}$



# History of Time Series



- early 1970s: ARIMA
- late 1970s: Bayesian Approach
- 1980s: GARCH
- 1995: Random Forest
- 1997: Support Vector Regression
- 1999: Boosted Decision Trees
- 2014: LSTMs, GRUs, Seq2Seq
- 2016: Dilated CNNs
- **2017: Transformer**

## Limitations of Statistical Methods

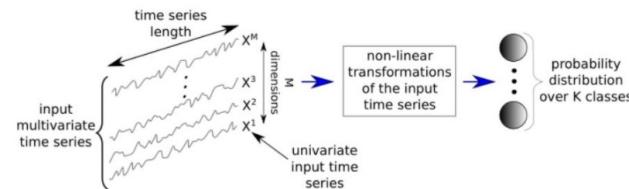
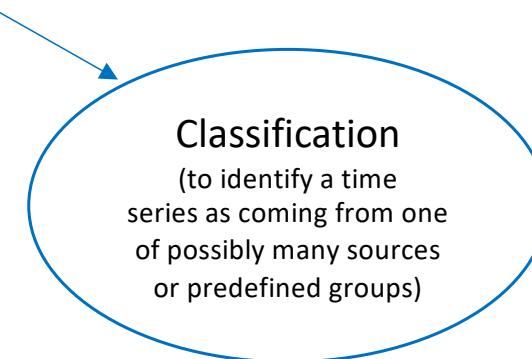
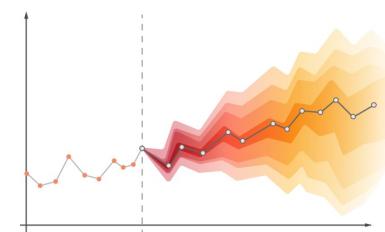
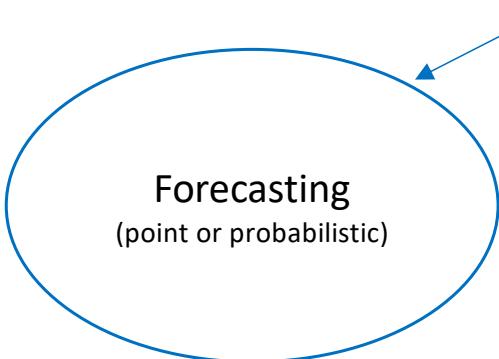
Focus on

- complete data
- linear relationships
- fixed temporal dependencies
- univariate data
- 1-step forecasts

## Why use Deep Learning for Time Series?

- Learns patterns in data with complex non-linear dependencies
- Supports multiple inputs and outputs
- Has shown good performance in many scenarios
- Models are flexible and expressive
- Can be trained with large datasets
- Easy to introduce exogenous features into the model

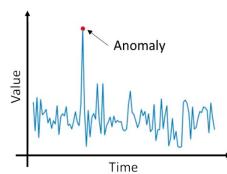
# Different tasks in Time series



## Other tasks:

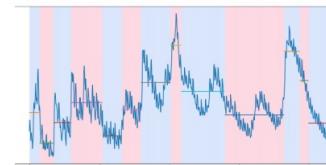
### 1. Anomaly Detection

*Find abnormal events in a time series*



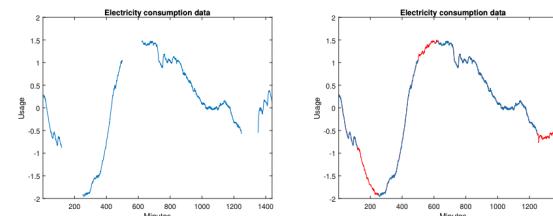
### 2. Segmentation/ change-point detection

*Find significant abrupt changes in the time series*



### 3. Completion/interpolation

*Recover missing/lost samples in a time series*



### 4. Query by content/indexation

*Given an input time series, retrieve the closest time series in a large database up to a given measure of fit*

# DL Architectures for Time Series Forecasting

## 1-step-ahead forecasting:

$$\hat{y}_{i,t+1} = f(y_{i,t-k:t}, \mathbf{x}_{i,t-k:t}, \mathbf{s}_i)$$


Model forecast      observations of target over a look-back window k      exogenous inputs over a look-back window k      static metadata associated with the entity (e.g. sensor location)

## multi-horizon forecasting:

$$\hat{y}_{t+\tau} = f(y_{t-k:t}, \mathbf{x}_{t-k:t}, \mathbf{u}_{t-k:t+\tau}, \mathbf{s}, \tau)$$

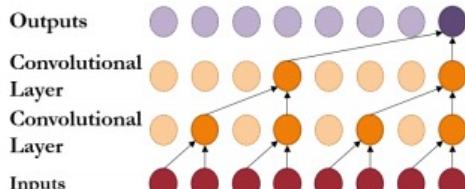


$\tau$  is a discrete forecast horizon

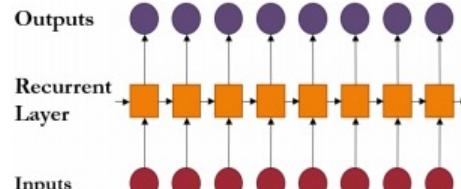
known future inputs (e.g. date information, such as the day-of-week or month) across the entire horizon

Time Series Forecasting With Deep Learning: A Survey, Bryan Lim and Stefan Zohren, <https://arxiv.org/pdf/2004.13408.pdf>

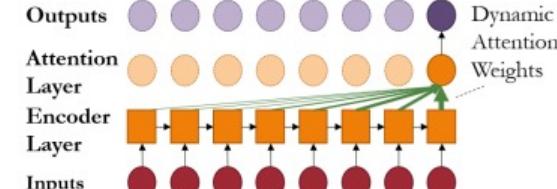
# Basic Building Blocks



(a) CNN Model.

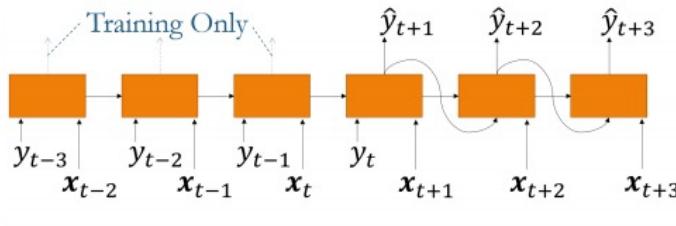


(b) RNN Model.



(c) Attention-based Model.

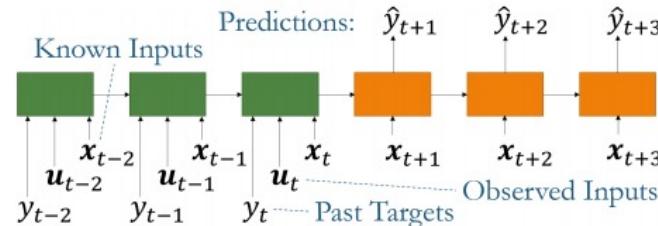
## Multi-horizon Forecasting Models



(a) Iterative Methods

autoregressive deep learning architectures, produce multi-horizon forecasts by recursively feeding samples of the target into future time steps

→ large error accumulations



(b) Direct Methods

produce forecasts directly using all available inputs, make use of sequence-to-sequence architectures, (an encoder summarizes past information and a decoder combines them with known future inputs)

→ Fixed maximum forecast horizon

Time Series Forecasting With Deep Learning: A Survey, Bryan Lim and Stefan Zohren, <https://arxiv.org/pdf/2004.13408.pdf>

# Outputs and Loss functions

1. **Point Estimates:** determine the expected value of a future target. This reformulates the problem to a classification task for discrete outputs (e.g. forecasting future events), and regression task for continuous outputs.

$$\mathcal{L}_{classification} = -\frac{1}{T} \sum_{t=1}^T y_t \log(\hat{y}_t) + (1 - y_t) \log(1 - \hat{y}_t)$$

$$\mathcal{L}_{regression} = \frac{1}{T} \sum_{t=1}^T (y_t - \hat{y}_t)^2$$

2. **Probabilistic Outputs:** model uncertainties and use deep neural networks to generate parameters of known distributions. e.g. Gaussian distributions are used to forecast continuous targets, with the networks outputting means and variance parameters for the predictive distributions at each step:

$$\mu(t, \tau) = \mathbf{W}_\mu \mathbf{h}_t^L + \mathbf{b}_\mu,$$

Final layer of the network

$$\zeta(t, \tau) = \text{softplus}(\mathbf{W}_\Sigma \mathbf{h}_t^L + \mathbf{b}_\Sigma),$$

Time Series Forecasting With Deep Learning: A Survey, Bryan Lim and Stefan Zohren, <https://arxiv.org/pdf/2004.13408.pdf>

# Convolutional Neural Networks

- Convolution layer
  - units in the hidden layer operate on a field of the output
  - weights are shared across input
- Motivation for CNNs in sequence learning
  - Success of CNN in CV and recently in NLP
  - Achieves state-of-the-art accuracy in audio synthesis and machine translation
- Main advantages:
  - Lower level of model complexity than RNNs
  - Parallelization of computations
  - Dilated CNN can outperform RNN in sequence modelling
- 2D Convolutions encode spatial invariance
- 1D Convolutions encode temporal invariance, “stationarity”

# 1D Convolutions

$$\begin{array}{c} \boxed{1 \ 3 \ 3 \ 0 \ 1 \ 2} \\ \underbrace{\hspace{1cm}}_{\text{Input } [1 \times 6]} \end{array} \times \begin{array}{c} \boxed{1 \ 0 \ 1} \\ \underbrace{\hspace{1cm}}_{\text{Filter } [1 \times 3]} \end{array} = \begin{array}{c} \boxed{4 \ 3 \ 4 \ 2} \\ \underbrace{\hspace{1cm}}_{\text{Output } [1 \times 4]} \end{array}$$

1D Convolution with 1 filter ( $stride = 1$ ),  $t = 6$ ,  $\#features = 1$

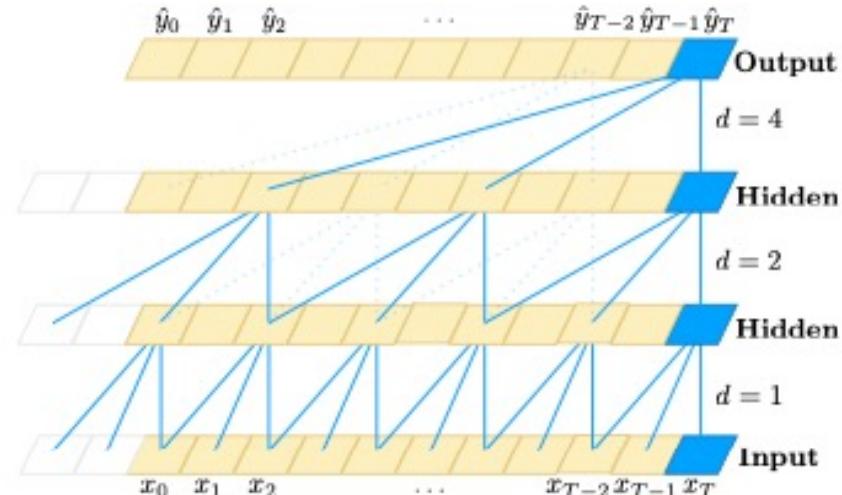
- A convolution kernel that is convolved with the input over a single temporal dimension
- Filters are trained to detect features in the sequence, regardless of where they appear
- Multiple filters can detect multiple features
- Can be implemented to multivariate time series (Input  $[t \times \#features]$ )

# Causality

- *causal constraint*:  $y_t$  depends only on  $x_0, \dots, x_t$  and not on any “future” inputs  $x_{t+1}, \dots, x_T$ .
- Learning: network  $f$  that minimizes some expected loss between the actual outputs and the predictions,  $L(y_0, \dots, y_T, f(x_0, \dots, x_T))$
- This formalism encompasses settings
  - auto-regressive prediction: predict some signal given its past:
    - setting the target output to be simply the input shifted by one time step.
  - does not capture domains (machine translation, or sequence-to-sequence prediction)
    - the entire input sequence (including “future” states) can be used to predict each output

# Causal Convolution

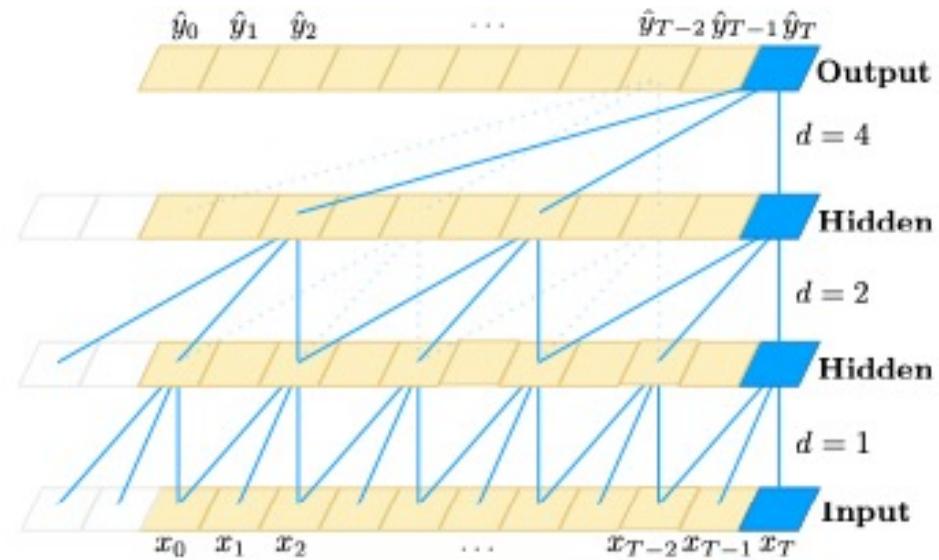
- Temporal Convolution Networks (TCN) constraints:
  - network produces an output of the same length as the input
    - To accomplish this: uses a 1D fully-convolutional network (FCN) architecture
    - each hidden layer is the same length as the input layer, and zero padding of length (kernel size – 1) is added to keep subsequent layers the same length as previous ones.
  - there can be no leakage from the future into the past.
    - TCN uses *causal convolutions*: output at time  $t$  is convolved **only** with elements  $\leq t$  in the previous layer.
  - $TCN = 1D\ FCN + causal\ convolutions$



An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling, Shaojie Bai. J. Zico Kolter, Vladlen Koltun, <https://arxiv.org/pdf/1803.01271.pdf>

# Causal Convolution

- Disadvantage: in order to achieve a long effective history size need large kernel size - extremely deep network or very large filters
- *Dilated convolution*: to allow for both very deep networks and very long effective history



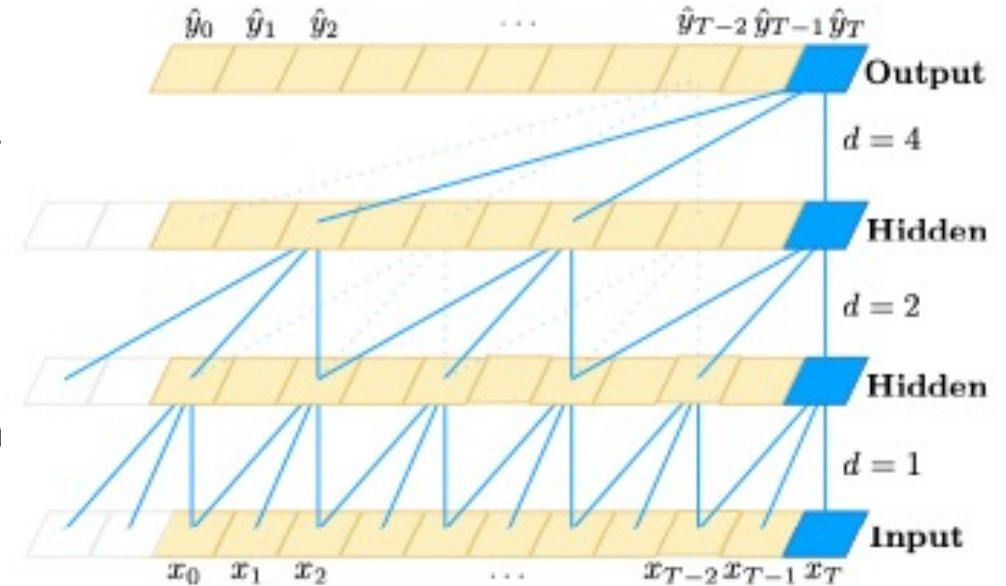
An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling, Shaojie Bai. J. Zico Kolter, Vladlen Koltun, <https://arxiv.org/pdf/1803.01271.pdf>

# Dilated Convolutions

- 1-D sequence input  $x \in R^n$  and a filter  $f:\{0,\dots,k-1\} \rightarrow R$ ,
- Dilated convolution operation  $F$  on element  $s$  of the sequence:

$$F(s) = (\mathbf{x} *_d f)(s) = \sum_{i=0}^{k-1} f(i) \cdot \mathbf{x}_{s-d \cdot i}$$

- $d$ : dilation factor,  $k$ : filter size,  $(s - d \cdot i)$  direction of the past.
- is equivalent to introducing a fixed step between every two adjacent filter taps.
- $d = 1$ : regular convolution
- $d > 1$ : enables output at the top level to represent a wider range of inputs: expanding the receptive field of a ConvNet.



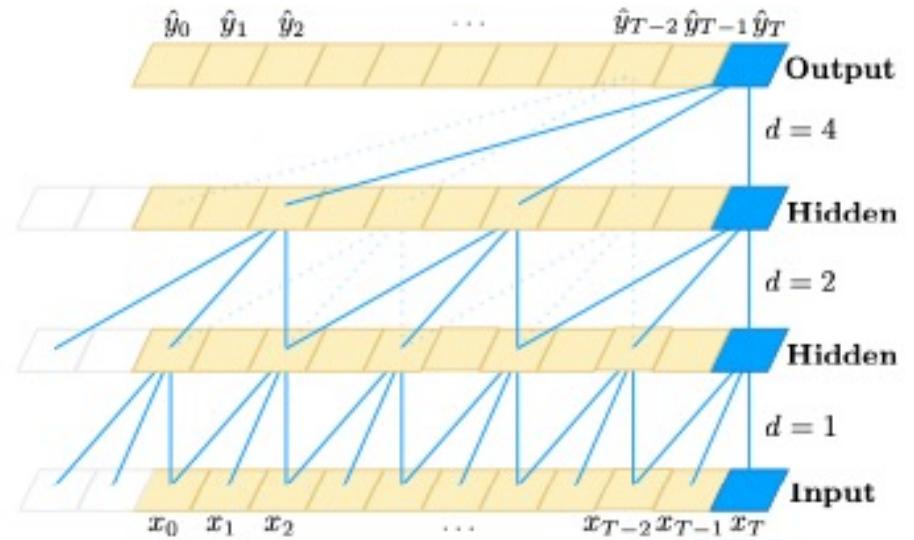
# Dilated Convolutions

Longer memory - increase the receptive field of the TCN:

- choosing larger filter sizes  $k$
- increasing the dilation factor  $d$ :
- effective history for this layer  $(k - 1)d$ .

In dilated convolutions, increase  $d$  exponentially with the depth of the network:  $d = O(2^i)$ .  $i$ : level in the network. Ensures

- each input treated by some filter in the effective history,
- allowing for an extremely large effective history using deep networks.



An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling,  
Shaojie Bai, J. Zico Kolter, Vladlen Koltun, <https://arxiv.org/pdf/1803.01271.pdf>

# Dilated Convolutions

Advantages TCNs for sequence modelling.

- *Parallelism*. Unlike in RNNs, convolutions can be done in parallel since the same filter is used in each layer.
- *Flexible receptive field size*: stacking more dilated (causal) convolutional layers, using larger dilation factors, increasing the filter size. better control of the model's memory size, and are easy to adapt to different domains.
- *Stable gradients*. Unlike recurrent architectures, TCN has a backpropagation path different from the temporal direction of the sequence. Thus avoids the problem of exploding/vanishing gradients
- *Low memory requirement for training*. For long input sequence, LSTMs and GRUs use a lot of memory for partial results for their multiple cell gates. TCN filters are shared across a layer, with the backpropagation path depending only on network depth.
- *Variable length inputs*. like RNNs, TCNs take in inputs of arbitrary lengths by sliding the 1D convolutional kernels.

Disadvantages

- *Data storage during evaluation*, RNNs only need hidden state  $h_t$  and the current input  $x_t$  in order to generate a prediction. In contrast, TCNs need the raw sequence up to the effective history length
- Potential parameter change for a transfer of domain.

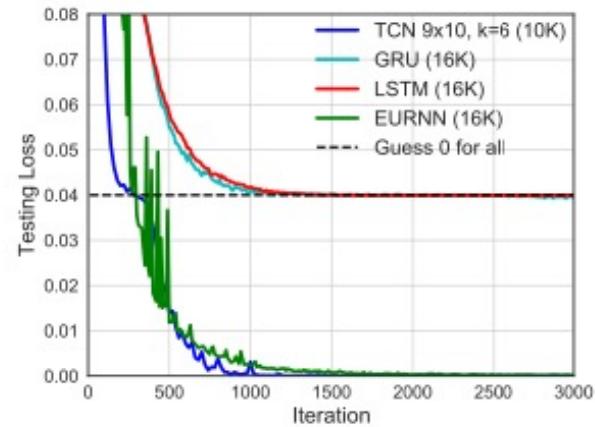
An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling, Shaojie Bai. J. Zico Kolter, Vladlen Koltun, <https://arxiv.org/pdf/1803.01271.pdf>

# Experimental Evaluation

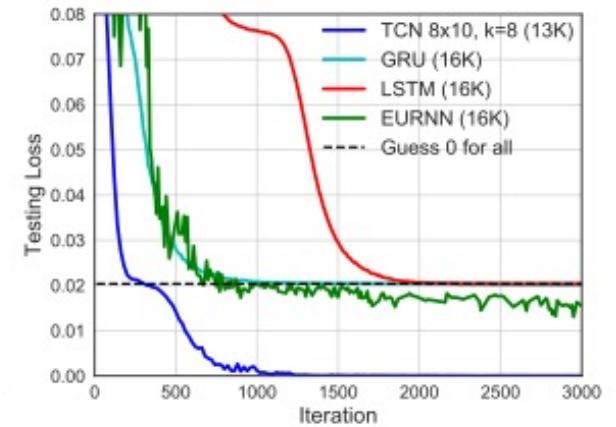
Sequence Modeling Task	Model Size ( $\approx$ )	Models			
		LSTM	GRU	RNN	TCN
Seq. MNIST (accuracy <sup>h</sup> )	70K	87.2	96.2	21.5	<b>99.0</b>
Permuted MNIST (accuracy)	70K	85.7	87.3	25.3	<b>97.2</b>
Adding problem $T=600$ (loss <sup>ℓ</sup> )	70K	0.164	<b>5.3e-5</b>	0.177	<b>5.8e-5</b>
Copy memory $T=1000$ (loss)	16K	0.0204	0.0197	0.0202	<b>3.5e-5</b>
Music JSB Chorales (loss)	300K	8.45	8.43	8.91	<b>8.10</b>
Music Nottingham (loss)	1M	3.29	3.46	4.05	<b>3.07</b>
Word-level PTB (perplexity <sup>ℓ</sup> )	13M	<b>78.93</b>	92.48	114.50	88.68
Word-level Wiki-103 (perplexity)	-	48.4	-	-	<b>45.19</b>
Word-level LAMBADA (perplexity)	-	4186	-	14725	<b>1279</b>
Char-level PTB (bpc <sup>ℓ</sup> )	3M	1.36	1.37	1.48	<b>1.31</b>
Char-level text8 (bpc)	5M	1.50	1.53	1.69	<b>1.45</b>

# Experimental Evaluation

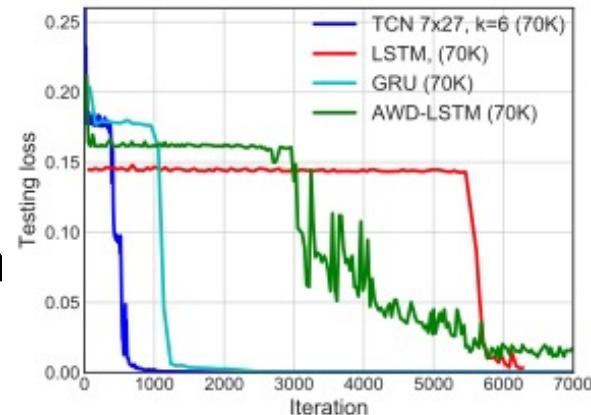
- Copy memory task



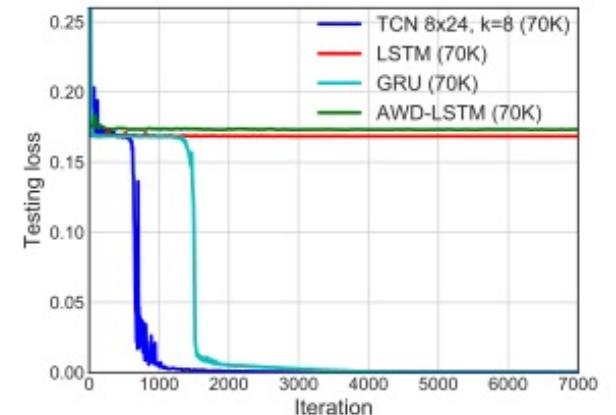
(a)  $T = 500$



(b)  $T = 1000$



(a)  $T = 200$

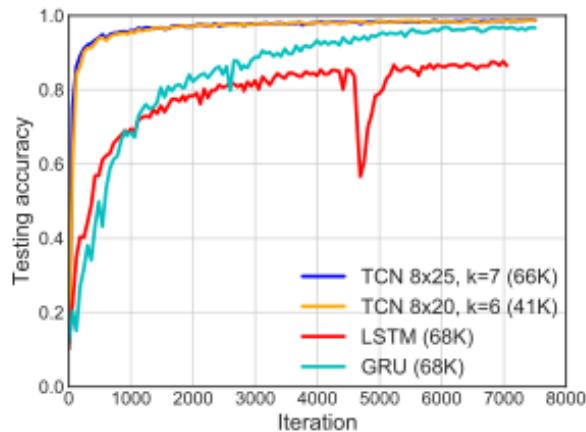


(b)  $T = 600$

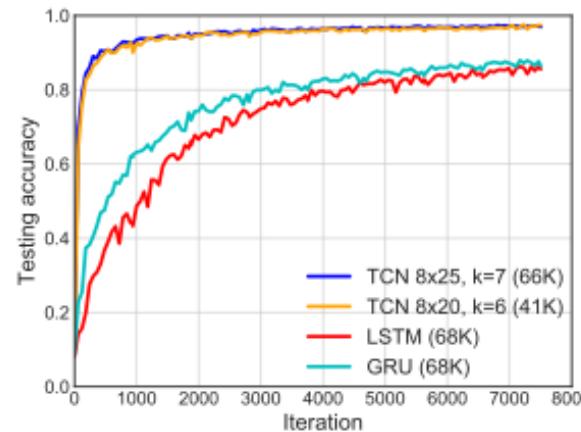
- adding problem task  
different sequence length

An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling, Shaojie Bai. J. Zico Kolter, Vladlen Koltun, <https://arxiv.org/pdf/1803.01271.pdf>

# Experimental Evaluation

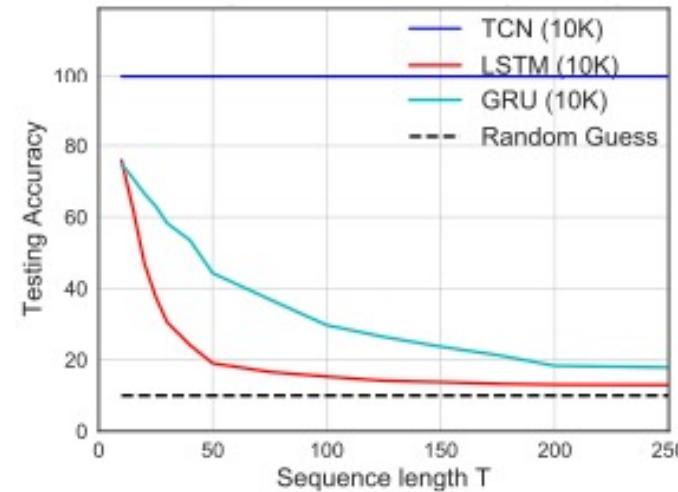


(a) Sequential MNIST



(b) P-MNIST

- Accuracy (copy memory task) for sequences of different lengths  $T$



# References

- [25 years of time series forecasting](#), Jan G. De Gooijer, Rob J. Hyndman b, 2006
- Deep Learning for Time-Series Analysis, 2017: <https://arxiv.org/abs/1701.01887>
- Deep Learning for Time-Series Forecasting: Predict the Future with MLPs, CNNs and LSTMs in Python, Jason Brownlee, Machine Learning Mastery, 2018
- FORECASTING ECONOMIC AND FINANCIAL TIME SERIES: ARIMA VS. LSTM, SIMA SIAMI NAMIN, <https://arxiv.org/pdf/1803.06386.pdf>, 2018.
- Deep Learning for Time-Series Analysis, Gamboa J., 2017, <https://arxiv.org/pdf/1701.01887.pdf>
- G.Peter Zhang, B.Eddy Patuwo, Michael Y. Hu, “A simulation study of artificial neural networks for nonlinear time-series forecasting,” Computers & Operations Research, Volume 28, Issue 4, 2001, pp. 381-396, ISSN 0305-0548, [https://doi.org/10.1016/S0305-0548\(99\)00123-9](https://doi.org/10.1016/S0305-0548(99)00123-9).
- G.Peter Zhang, Min Qi, “Neural network forecasting for seasonal and trend time series,” European Journal of Operational Research, Volume 160, Issue 2, 2005, pp. 501-514, ISSN 0377-2217, <https://doi.org/10.1016/j.ejor.2003.08.037>.
- <https://machinelearningmastery.com/remove-trends-seasonality-difference-transform-python/>
- Long, Jonathan, Shelhamer, Evan, and Darrell, Trevor. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015
- Bryan Lim and Stefan Zohren, Time Series Forecasting With Deep Learning: A Survey, <https://arxiv.org/abs/2004.13408>

# References

- MLP
  - Shiblee, M., Kalra, P. K., & Chandra, B. (2009). Time Series Prediction with Multilayer Perceptron (MLP): A New Generalized Error Based Approach. Lecture Notes in Computer Science, 37–44. doi:10.1007/978-3-642-03040-6\_5
  - N-BEATS: Neural basis expansion analysis for interpretable time series forecasting, Boris N. Oreshkin, Dmitri Carov, Nicolas Chapados and Yoshua Bengio, <https://arxiv.org/pdf/1905.10437.pdf>
- RNNs, LSTMs
  - K. A. Althelaya, E. M. El-Alfy and S. Mohammed, "Evaluation of bidirectional LSTM for short-and long-term stock market prediction," 2018 9th International Conference on Information and Communication Systems (ICICS), Irbid, 2018, pp. 151-156, doi: 10.1109/ICACS.2018.8355458.
  - DeepAR: Probabilistic Forecasting with Autoregressive Recurrent Networks, David Salinas, Valentin Flunkert and Jan Gasthaus, <https://arxiv.org/pdf/1704.04110.pdf>
  - Deep State Space Models for Time Series Forecasting, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi and R. Garnett, <https://papers.nips.cc/paper/2018/file/5cf68969fb67aa6082363a6d4e6468e2-Paper.pdfdeep>
- 1D CNN, Causal CNN, Dilated CNN
  - \*Oord et al. WaveNet: A Generative Model for Raw Audio. <https://arxiv.org/pdf/1609.03499.pdf>
  - \*S. Bai, J. Kolter, and V. Koltun. An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling. <https://arxiv.org/pdf/1803.01271.pdf>
  - Chang et al. Dilated Recurrent Neural Networks.<https://papers.nips.cc/paper/2017/file/32bb90e8976aab5298d5da10fe66f21d-Paper.pdf>
- Hybrid Models (ARIMA-LSTM, LSTM-CNN, Attention-LSTM)
  - Modeling Long- and Short-Term Temporal Patterns with Deep Neural Networks, Wuokun et al, 2018, <https://dl.acm.org/doi/pdf/10.1145/3209978.3210006>
  - A Dual-Stage Attention-Based Recurrent Neural Network for Time Series Prediction, Yao Qin, Dongjin Song, Haifeng Chen, Wei Cheng, Guofei Jiang and Garrison Cottrell, <https://songdj.github.io/publication/ijcai-17-a/ijcai-17-a.pdf>
  - Think Globally, Act Locally: A Deep Neural Network Approach to High-Dimensional Time Series Forecasting, Rajat Sen, Hsiang-Fu Yu and Inderjit Dhillon, <https://papers.nips.cc/paper/2019/file/3a0844cee4fcf57de0c71e9ad3035478-Paper.pdf>
  - Hybrid Neural Networks for Learning the Trend in Time Series, Lin et al. 2017, <https://infoscience.epfl.ch/record/262447>
  - Learning representations from eeg with deep recurrent-convolutional neural networks, Bashivan et al., arXiv preprint arXiv:1511.06448, 2015, <https://arxiv.org/pdf/1511.06448.pdf>
  - Wang, R., Peng, C., Gao, J. et al. A dilated convolution network-based LSTM model for multi-step prediction of chaotic time series. Comp. Appl. Math. 39, 30 (2020). <https://doi.org/10.1007/s40314-019-1006-2>
  - G.Peter Zhang. "Time series forecasting using a hybrid ARIMA and neural network model", Neurocomputing, Volume 50, 2003, pp. 159-175, ISSN 0925-2312, [https://doi.org/10.1016/S0925-2312\(01\)00702-0](https://doi.org/10.1016/S0925-2312(01)00702-0).
  - Bao W, Yue J, Rao Y (2017) A deep learning framework for financial time series using stacked autoencoders and long-short term memory. PLoS ONE 12(7): e0180944. <https://doi.org/10.1371/journal.pone.0180944>
- Transformer for Time-series
  - Enhancing the Locality and Breaking the Memory Bottleneck of Transformer on Time Series Forecasting, Li et al., 2019, <https://arxiv.org/pdf/1907.00235.pdf>
- Graphs and GNNs
  - Diffusion Convolutional Recurrent Neural Network: Data-Driven Traffic Forecasting, Yaguang Li, Rose Yu, Cyrus Shahabi and Yan Liu, <https://arxiv.org/pdf/1707.01926.pdf>
  - Spectral Temporal Graph Neural Network for Multivariate Time-series Forecasting, Defu Cao, Yujing Wang, Juanyong Duan, Ce Zhang, Xia Zhu, Conguri Huang, Yunhai Tong, Bixiong Xu, Jing Bai, Jie Tong and Qi Zhang, <https://arxiv.org/pdf/2103.07719.pdf>

# AI for Recommendations

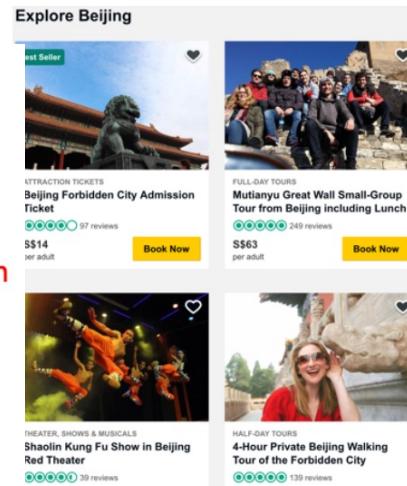
# AI for Recommendations

Recommendation widely applied online:

- E-commerce (Amazon)
- Content Sharing (Netflix, Youtube, Pinterest)
- Social Networking (Facebook, Linked-in Twitter)
- Forums (Reddit, IMDB,
- Gaming (Tencent)
- ....



Image & Video  
Recommendation



Screenshot of TripAdvisor

POI & Post  
Recommendation

The screenshot shows a search results page on Taobao for "iphone11". It includes:

- Product Recommendation:** A listing for a "隐形壳" (Invisible Case) for the iPhone 11 Pro Max, priced at ¥68, with a rating of 4.8 stars and over 650,000 reviews. It's marked as a "5-star old store" and "best seller".
- Product Recommendation:** A listing for the "Apple iPhone 11", priced at ¥6799, with a rating of 4.8 stars and over 240,000 reviews. It's marked as a "best seller".
- Friend Recommendation:** A section titled "People You May Know" showing profiles of users with 2, 3, 4, and 5 mutual friends.

Ad & Product  
Recommendation

Friend  
Recommendation

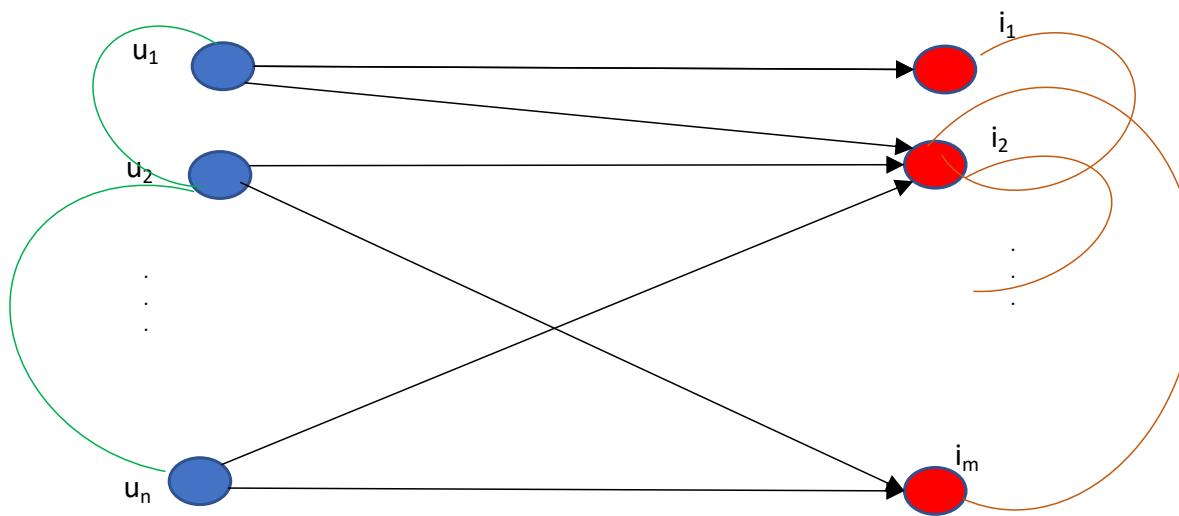
# Value of Recommendation Systems

- Netflix: 80% of what people watch comes from some sort of recommendation.
  - 2015: “*We think the combined effect of personalization and recommendations save us more than \$1B per year.*”  
*Netflix paper:* <https://dl.acm.org/doi/pdf/10.1145/2843948>
- Amazon: 30%+ page views comes from recoms[Smith and Linden, 2017]
- YouTube Homepage: 60%+ clicks [Davidson et al. 2010]
- Airbnb, Search Ranking ....: 99% of all booking conversions.

## Early days

- **Netflix**: personalized video-recommendation system based on ratings and reviews by its customers.
- In 2006, offered a \$1,000,000 prize to the first developer of a video-recommendation algorithm that could beat its existing algorithm - *Cinematch* by 10%.
- Collaborative filtering early Recommender Systems
  - Spotify
  - YouTube.
  - LinkedIn (Browsemaps).

## Basic Info: Users – Items graphs



### User-User

- Social Relations
- Same Profiles

...

### User-Item

- Implicit Feedback
- Explicit Feedback ...

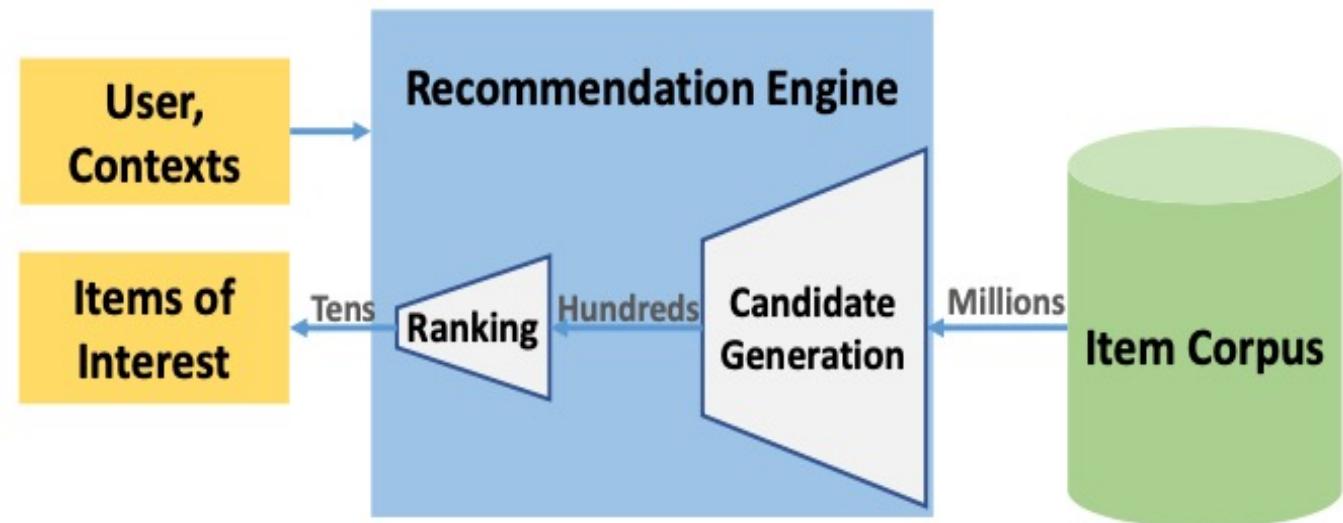
### Item-Item

- Same Attributes
- External Knowledge ...

# Recommendation Engine architecture

## User Contexts

- Interaction history
- Demographics ...



## Items

- Products, News, Movies,  
Videos, Friends ...

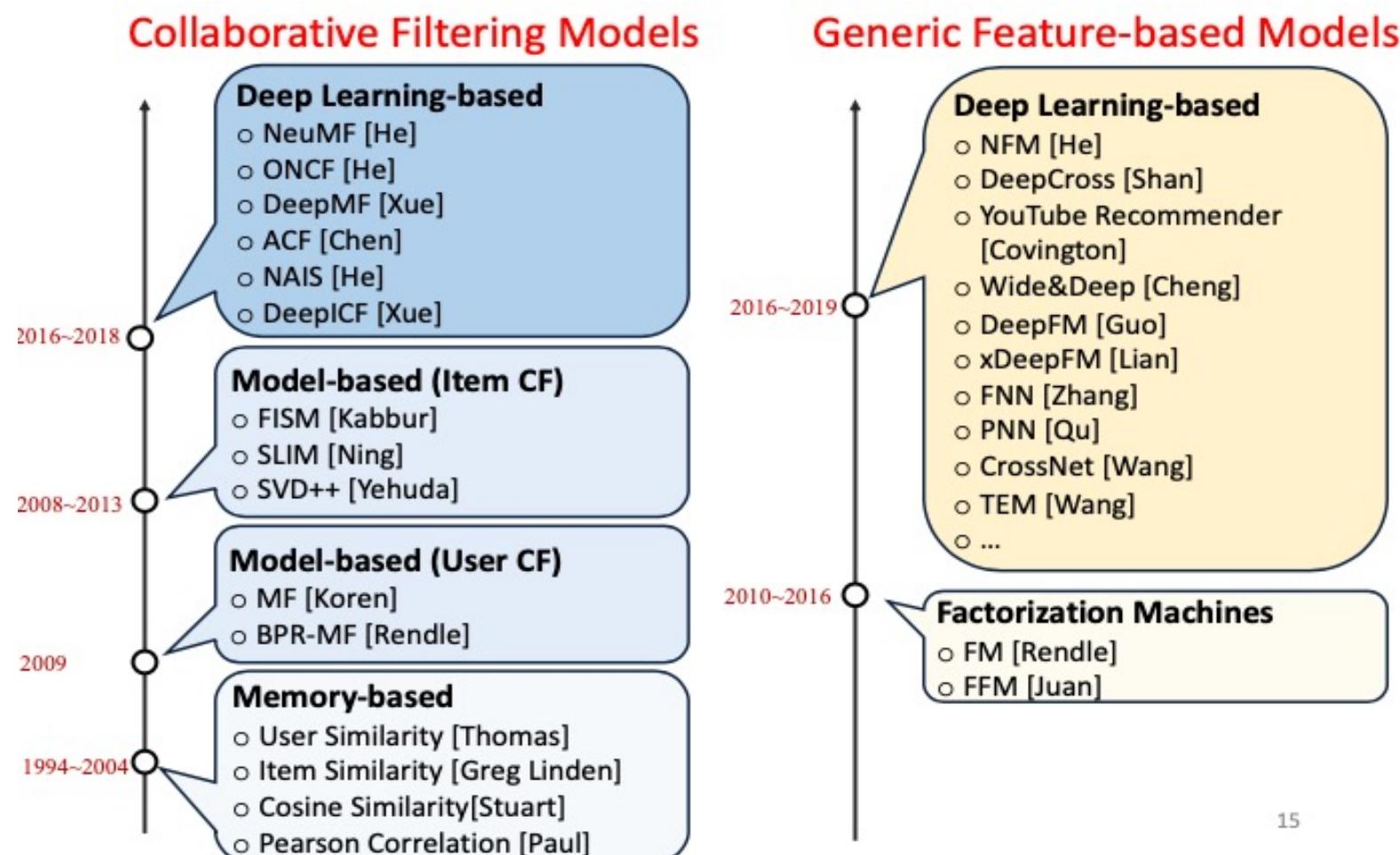
## Recommender Data Model

- Set  $U=\{u_1, \dots, u_n\}$  of users
- Set  $I=\{i_1, \dots, i_m\}$  of items (e.g. products)
- Elements from  $U$  and  $I$  can be described by a vector respectively
  - $(a_1, \dots, a_s) \rightarrow$  attributes of user profile
  - $(b_1, \dots, b_t) \rightarrow$  description of items (meta data, features, ...)
- Goal of recommendation process: recommend **new** items for an active user  $u$
- Overview of process
  - User modeling (explicit or implicit, e.g. user rates items)
  - Personalization, generate list of recommended items

# Issues of Recommender Systems

- Cold start and latency problems
- Sparseness of user-item matrix
- Diversity of recommendations
- Scalability
- Privacy and trust
- Robustness
- Utilization of domain knowledge
- Changing user interests (dynamics)
- Evaluation of recommender systems

# Recommendations methods landscape



15

Learning and Reasoning on Graph for Recommendation, Xiang Wang , Xiangnan He, Tat-Seng Chua , tutorial CIKM 2019

# Recommendations

## Introduction

### Methods - Background

- Collaborative filtering
- Matrix completion NMF

### Advanced Methods

- Factorization machines –
- Deep Learning

# Collaborative Filtering (CF)

- Basic idea: System recommends items which were preferred by **similar** users in the past
  - Based on ratings
    - Expressed preferences of the active user
    - And also other users → Collaborative approach
  - Works on user-item matrix
    - Memory-based or model-based
    - No item meta data etc.!
- Assumption: Similar taste in the past implies similar taste in future
- CF is formalization of “word of mouth” among buddies

# General Process

1. Users rate items
2. Find set  $S$  of users which have rated similar to the active user  $u$  in the past ( $\rightarrow$  neighborhood)
  - Similarity calculation
  - Select the  $k$  nearest users to the active user
3. Generate candidate items for recommendation
  - Items which were rated in neighborhood of  $u$ ,
  - but were not rated by  $u$  yet
4. Predict rating of  $u$  for candidate items
  - Select and display  $n$  best items

## Example (I)

	Hoop Dreams	Star Wars	Pretty Woman	Titanic	Blimp	Rocky XV
Joe	D	A	B	D	?	?
John	A	F	D		F	
Susan	A	A	A	A	A	A
Pat	D	A		C		
Jean	A	C	A	C		A
Ben	F	A				F
Nathan	D		A		A	

Source: <http://www.dfki.de/~jameson/ijcai03-tutorial/>

## Example (II)

	Hoop Dreams	Star Wars	Pretty Woman	Titanic	Blimp	Rocky XV
Joe	D	A	B	D	?	?
John	A	F	D		F	
Susan	A	A	A	A	A	A
Pat	D	A		C		
Jean	A	C	A	C		A
Ben	F	A				F
Nathan	D		A		A	

## Example (III)

	Hoop Dreams	Star Wars	Pretty Woman	Titanic	Blimp	Rocky XV
Joe	D	A	B	D	?	?
John	A	F	D		F	
Susan	A	A	A	A	A	A
Pat	D	A		C		
Jean	A	C	A	C		A
Ben	F	A				F
Nathan	D		A		A	

# Memory based CF

		item					
		1	2	3	4		
user	1	5	?	?	?	...	...
	2	3	4	?	?		
3	?	1	2	4	...	...	...
...	...	...	...	...	...	...	...

Interaction Matrix

Problem: predict user  $u$ 's rating on item  $i$ .

- User-based CF computes ratings of  $u$ 's *similar users* on the target item  $i$ .

$$\hat{y}_{ui} = \sum_{u' \in S_u(u)} sim(u, u') \cdot y_{u'i}$$

Rating of a similar user on  $i$

Similar users of  $u$

- Item-based CF leverages the ratings of  $u$  on other *similar items* of  $i$

$$\hat{y}_{ui} = \sum_{i' \in S_i(i)} sim(i, i') \cdot y_{ui'}$$

Rating of  $u$  on a similar item

Similar items of  $i$

- Many similarity measures can be used, e.g., Jaccard, Cosine, Pearson Correlation. Recent advance learns the similarity from data.

# Recommendations

## Introduction

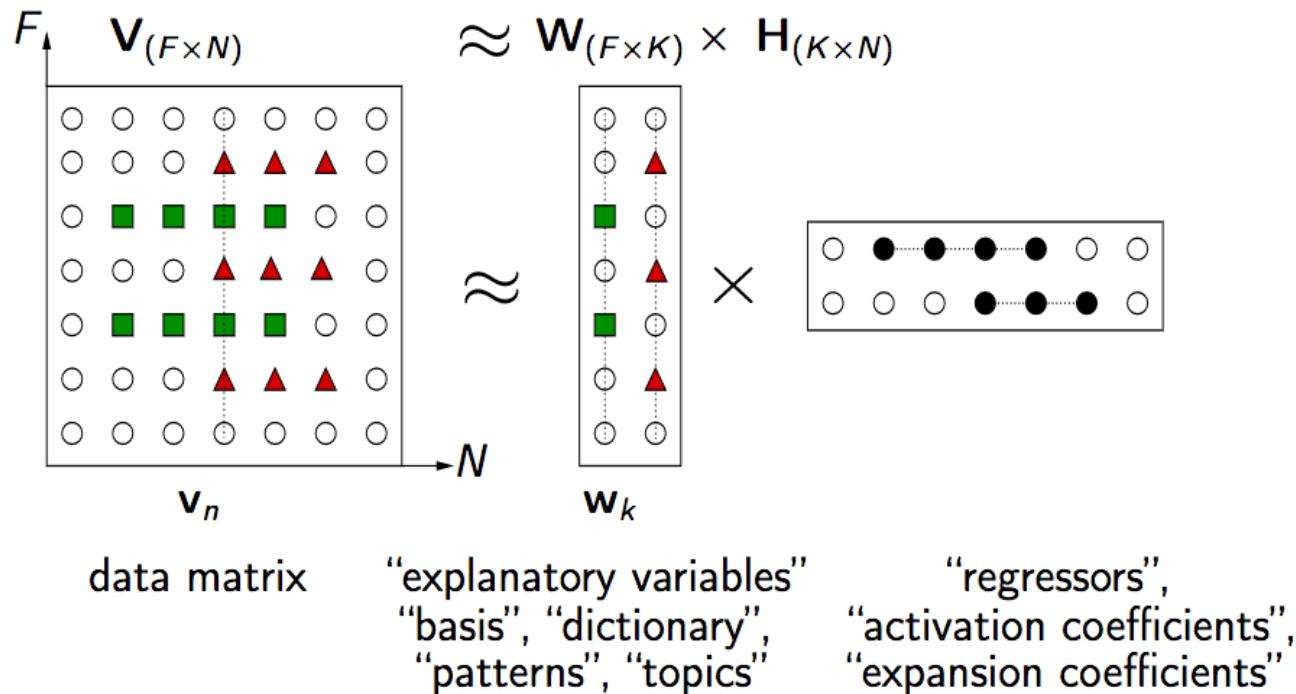
### Methods - Background

- Collaborative filtering
- **Matrix completion NMF**

### Advanced Methods

- Factorization machines –
- Deep Learning

# Explaining data by factorization



*Illustration by C. Févotte*

# Low Rank approximation

- Consider the rating matrix  $R$  with elements  $R_{ui}$   
If we have a set of users  $U$  and a set of items  $I$ ,
- Rating Matrix is **sparse**
- then  $R$  is a  $|U| \times |I|$  matrix
- We can approximate this by a low-rank approximation:
- $R \approx AB$
- where  $A$  is a  $|U| \times K$  matrix and  $B$  is a  $K \times |I|$  matrix
- $K$  is the rank of  $AB$  (number of non-zero eigenvalues)

# Low Rank approximation

- ratings  $R_{ui}$  for  $(u,i) \in D$  where  $|D| \ll |U| \times |I|$  (*number of ratings is typically much less than the total possible number of ratings*)
- In low rank approximation we are trying to fit to  $|D|$  ratings of the matrix with  $(|U| + |I|) \times K$  elements
- $K$  controls the complexity of the model
- Large  $K$  will fit the data better but at the cost of possibly over-fitting

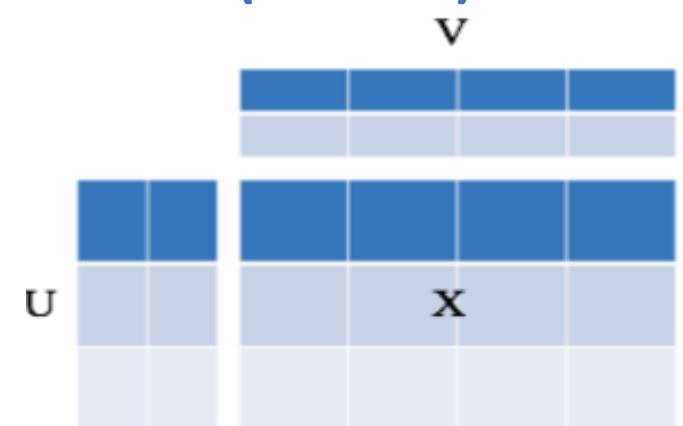
# Non Negative Matrix factorization (NMF)

- Data is often nonnegative by nature
  - pixel intensities; occurrence counts; food or energy consumption; user scores; stock market values;
- Interpretability of the results, optimal processing of nonnegative data may call for processing under Non-negativity constraints
- Applying SVD results in factorized matrices with positive and negative elements may contradict the physical meaning of the result.
  - - *Nonnegative matrix factorization (NMF)*
  - find the reduced rank *nonnegative factors* to approximate a given nonnegative data matrix.

# Non Negative Matrix factorization (NMF)

- Assume  $X$  ( $m \times n$ ) data matrix and  $r \ll m, n$
- NMF aims to find non negative matrices

$$U \in R^{m \times r}, V \in R^{r \times n} : X \approx UV^T$$



- To find  $U, V$ , optimization problem:

$$\min_{(U,V)} \|X - UV^T\|_2$$

$$X \simeq UV^T$$

- $U = [u_{fk}], w_{fk} \geq 0$
- $V = [v_{kn}], h_{kn} \geq 0$
- $k \ll f, n$

- Alternative error function:

$$\begin{aligned} \min_{U,V} f(U, V) &= \sum_{i=1}^m \sum_{j=1}^n \left( X_{ij} \log \frac{X_{ij}}{(UV^T)_{ij}} - X_{ij} + (UV^T)_{ij} \right) \\ \text{s.t. } U_{ia} &\geq 0, V_{jb} \geq 0, \forall i, a, b, j \end{aligned}$$

# NMF - Alternating Least squares

1. Suppose we know  $U$ , with  $V$  unknown.

- for each  $j$  we could minimize  $\|X_{\cdot j} - UV_{\cdot j}^T\|_2$ 
  - find  $V_{\cdot j}$  that minimizes with  $X_{\cdot j}$  and  $U$  known.
  - Frobenius norm: sum of squares,
  - minimization is a least-squares problem, i.e. linear regression
  - “predicting”  $X_{\cdot j}$  from  $W$ .

$$V_{\cdot j} = (U^T U)^{-1} U^T X_{\cdot j}$$

- repeat for all columns  $V_{\cdot j}$

2. assume  $V$ , with  $U$  unknown:  $X^T = VU^T$

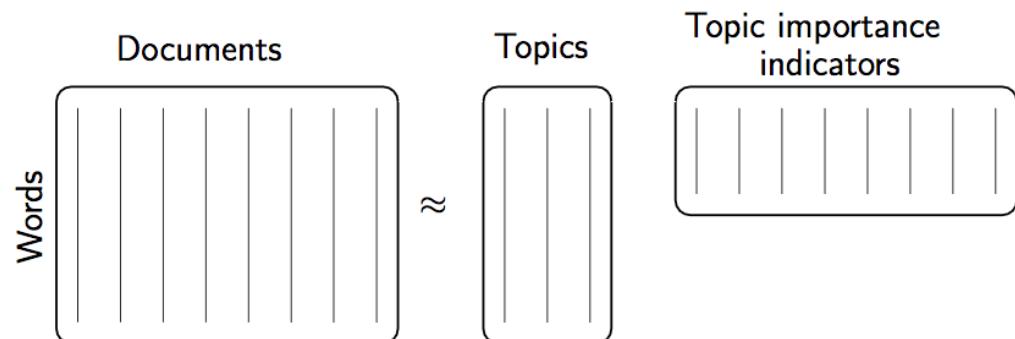
- Interchange roles of  $U$ ,  $V$  in the above optimization
- Compute a row of  $U$ , repeat for all rows

# NMF - Alternating Least squares

- Putting all this together
  - first choose initial guesses, random numbers, for U and V
  - alternate:
    - Compute U assuming V known
    - Compute V based on that new U
    - ...
- may generate some negative values: simply truncate to 0

# NMF issues, applications

- Uniqueness and Convergence
- $U_{mxr}$ ,  $r$  (rank) choice: via SVD...
- Applications
  - Topic detection
  - Source separation (music , speech)
  - Clustering
  - **Recommendations**



# Recommendations

## Introduction

### Methods - Background

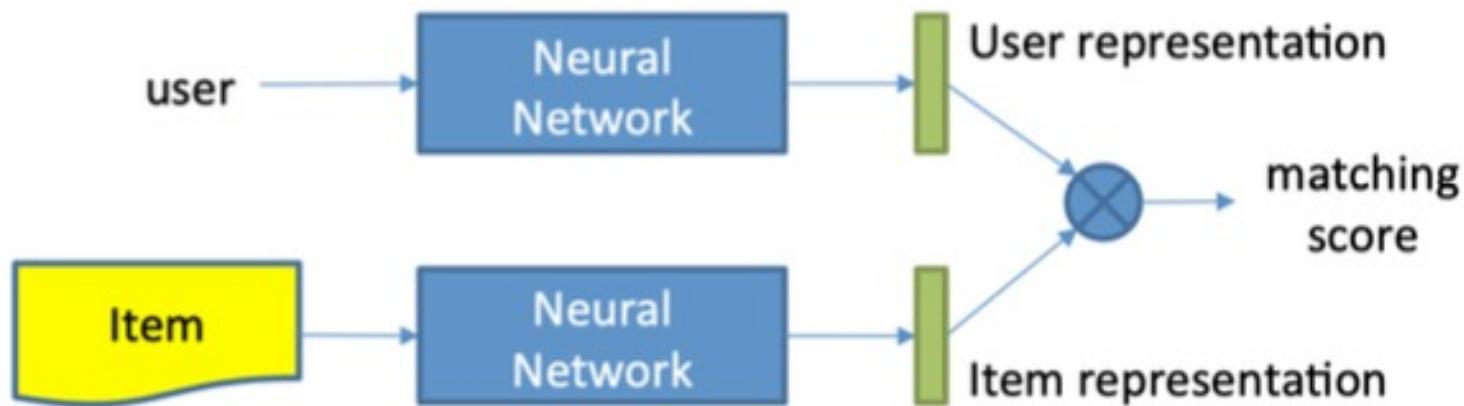
- Collaborative filtering
- Matrix completion NMF

### Advanced Methods

- Deep Learning

# Deep Learning for CF

- deep learning to embed users, item and their interaction in an embedding space

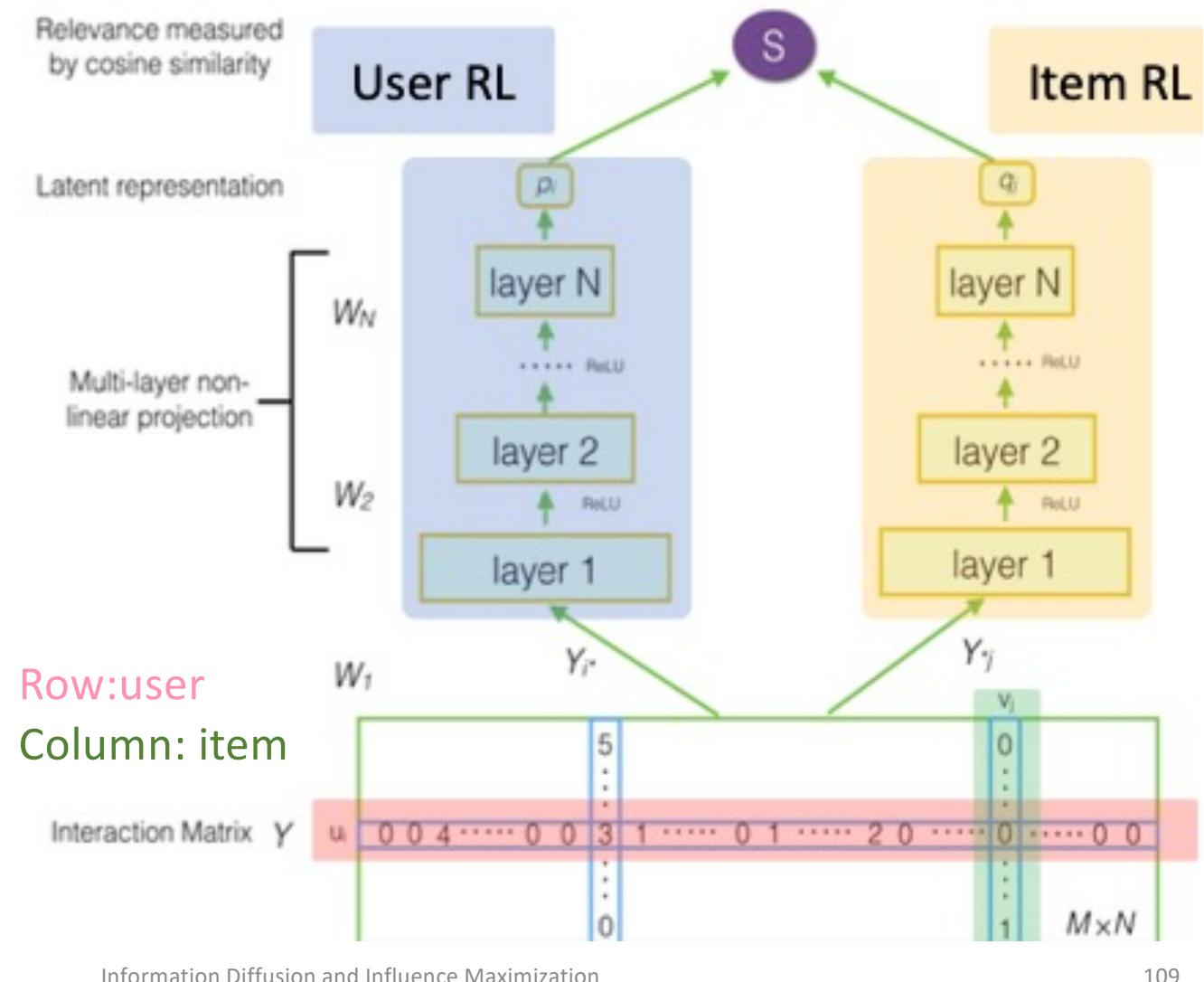


Model	Input Data	Representation Learning	Interaction Learning
DeepMF [Xue, IJCAI'17]	User: Historical items Item: User group	Multi-Layer Perceptron	Inner product
AutoRec [Sedhain, WWW'15]	User: Historical items Item: ID	Multi-Layer Perceptron	Inner Product
CDAE [Wu, WSDM'16]	User: Historical items + ID Item: ID	Multi-Layer Perceptron	Inner Product

23

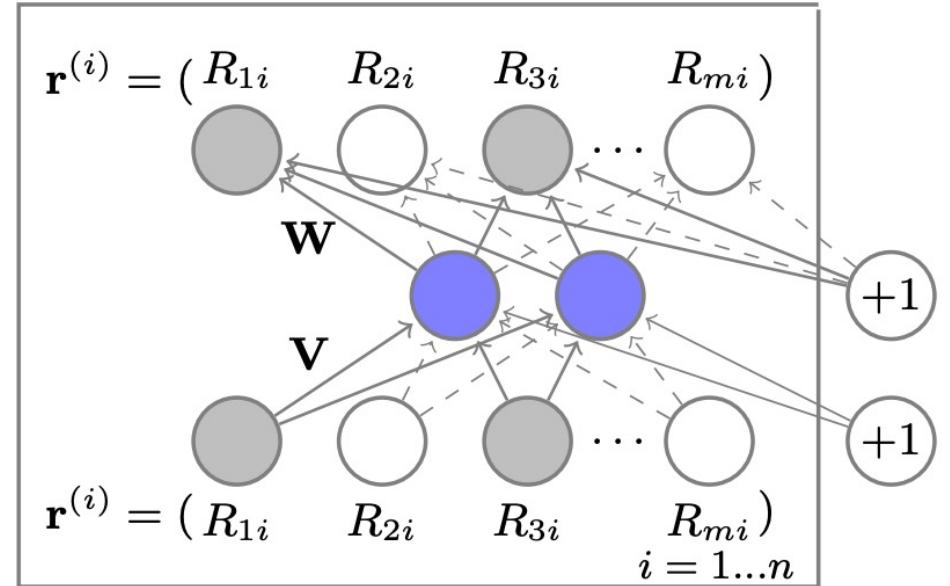
# Deep Matrix Factorization (Xue, IJCAI'17)

- Representation Learning Multi-layer perceptron
  - Deep Neural Networks are adopted to learn representations of
- users & items



# AutoRec: Autoencoders Meet Collaborative Filtering

- $m$  users,  $n$  items,
- $R \in R^{m \times n}$ : partially observed user-item rating matrix .
- user  $u \in U = \{1 \dots m\}$ :  $r(u) = (R^{ux1}, \dots, R^{un}) \in R^n$  .
- item  $i \in I = \{1 \dots n\}$   $r(i) = (R^{1i}, \dots, R^{mi}) \in R^m$ .
- item-based (user-based) autoencoder
- input each partially observed  $r(i)$  ( $r(u)$  ),
- project it into a low-dimensional latent (hidden) space
- reconstruct  $r(i)$  ( $r(u)$ ) in the output space to **predict missing ratings** for purposes of recommendation



$$\min_{\theta} \sum_{\mathbf{r} \in \mathbf{S}} \|\mathbf{r} - h(\mathbf{r}; \theta)\|_2^2.$$

$h(r; \theta)$  is the reconstruction of input  $r \in R^d$

$$h(\mathbf{r}; \theta) = f(\mathbf{W} \cdot g(\mathbf{V}\mathbf{r} + \boldsymbol{\mu}) + \mathbf{b})$$

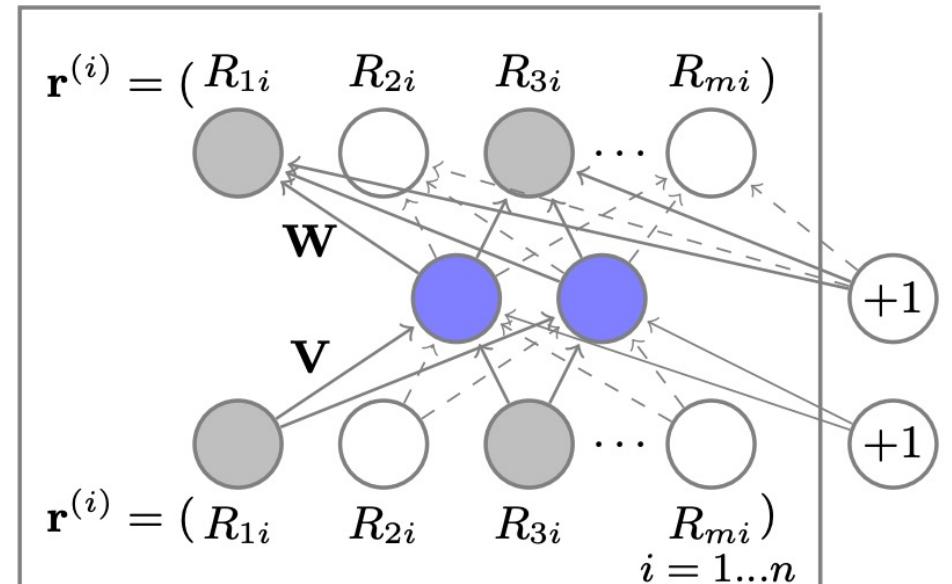
# Auto rec

- the objective function for the Item-based AutoRec (I-AutoRec) model is, for regularisation  $\lambda > 0$

$$\min_{\theta} \sum_{i=1}^n \|\mathbf{r}^{(i)} - h(\mathbf{r}^{(i)}; \theta)\|_{\mathcal{O}}^2 + \frac{\lambda}{2} \cdot (\|\mathbf{W}\|_F^2 + \|\mathbf{V}\|_F^2).$$

- predicted rating for user  $u$  and item  $i$

$$\hat{R}_{ui} = (h(\mathbf{r}^{(i)}; \hat{\theta}))_u$$



# Auto rec - Experiments

	ML-1M	ML-10M	$f(\cdot)$	$g(\cdot)$	RMSE
U-RBM	0.881	0.823	Identity	Identity	0.872
I-RBM	0.854	0.825	Sigmoid	Identity	0.852
U-AutoRec	0.874	0.867	Identity	Sigmoid	<b>0.831</b>
I-AutoRec	<b>0.831</b>	<b>0.782</b>	Sigmoid	Sigmoid	0.836

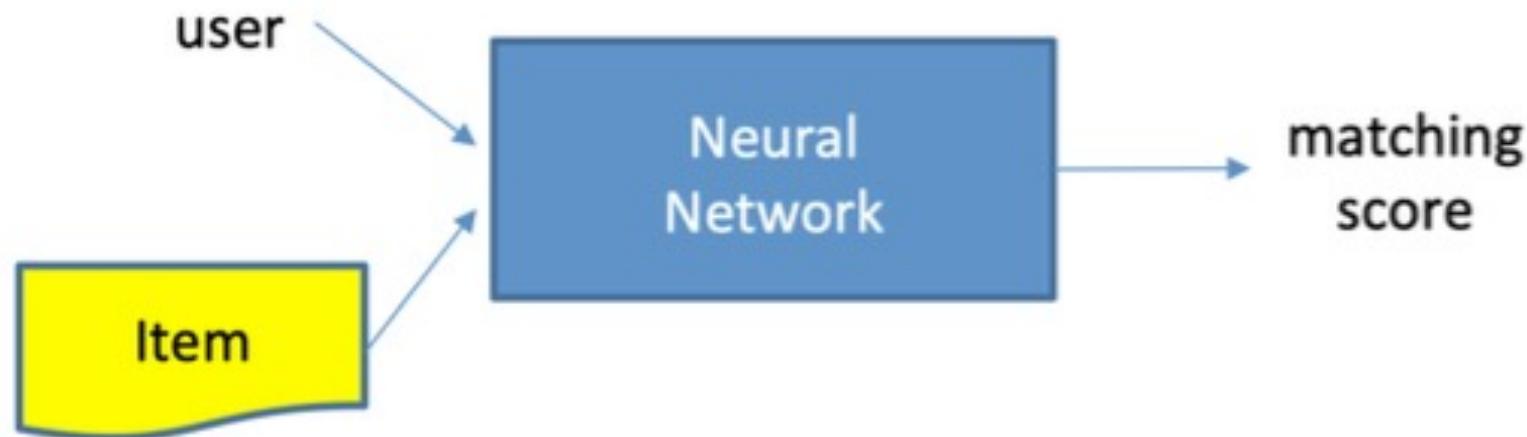
- RMSE of I/U-AutoRec and RBM models.
- RMSE for I-AutoRec different activation functions, MovieLens 1M dataset.

	ML-1M	ML-10M	Netflix
BiasedMF	0.845	0.803	0.844
I-RBM	0.854	0.825	-
U-RBM	0.881	0.823	0.845
LLORMA	0.833	<b>0.782</b>	0.834
I-AutoRec	<b>0.831</b>	<b>0.782</b>	<b>0.823</b>

- I-AutoRec vs baselines on MovieLens and Netflix datasets.

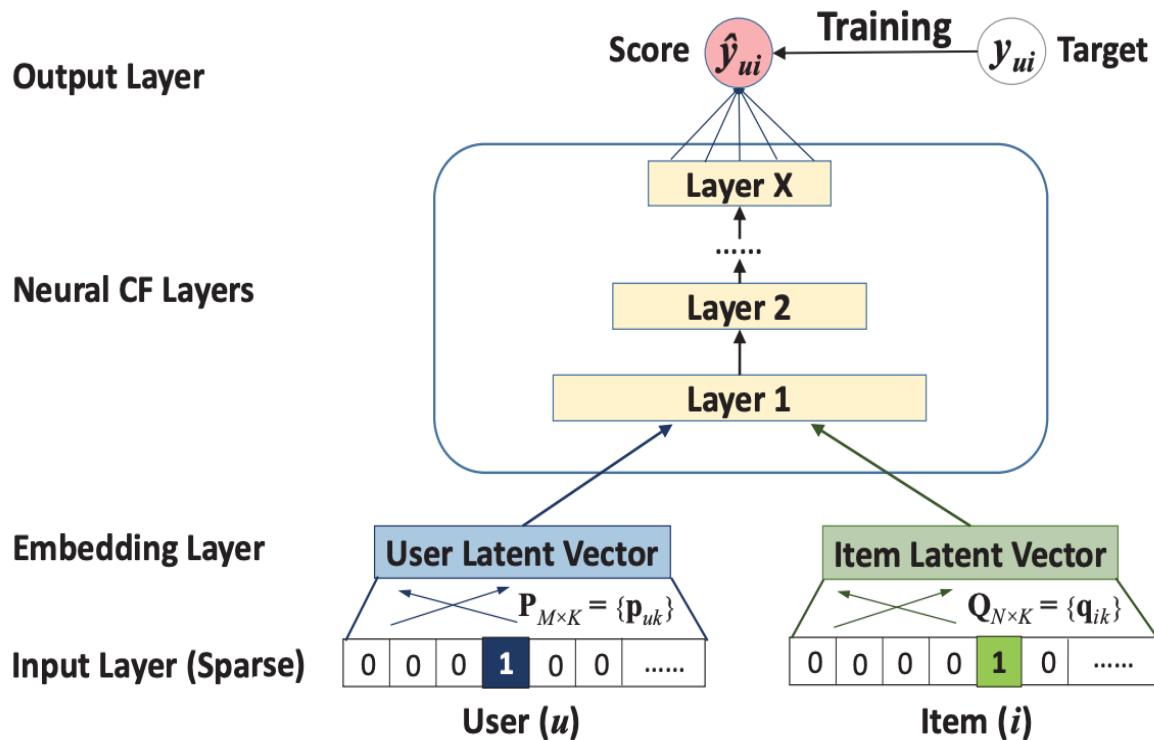
# Deep Learning for CF II

- Methods of interaction function learning
  - Capture complex patterns of user-item relationships



Model	Input Data	Representation Learning	Interaction Learning
NCF [He, WWW'17]	User: ID Item: ID	ID embedding	Multi-Layer Perceptron
NNCF [Bai, CIKM'17]	User: User neighbors Item: Item neighbors	Embeddings	Multi-Layer Perceptron
ONCF [He, IJCAI'28]	User: ID Item: ID	ID embedding	Convolutional Neural Network

# Neural Matrix Factorization (He, WWW'17)



## Interaction Modelling

- MF + MLP over users and items
- MF uses inner product to capture the low-rank relation
- DNN to learn the matching function

# Recommendation evaluation metrics

- top-K Hit Ratio (HR@K): measures whether the test item appears on Top-K list,
- Normalized Discounted Cumulative Gain (NDCG)  
Replaces hit with allocating higher scores to hits with top ranks

HR@K: a recall-based measure, as,

$$HR@K = \frac{\#hits@K}{|GT|}$$

NDCG: is a ranking-based measure, as,

$$NDCG@K = Z_k \sum_{i=1}^K \frac{2^{rel_i} - 1}{\log_2(i + 1)}$$

- $GT$ : the test list set,
- $rel_i$ : graded relevance value of the item at position  $i$ .  $rel_i \in \{0, 1\}$ , which depends on whether  $i$  is in the test dataset.
- $Z_K$ : normalization.

# NeurMF - Experiments

Dataset	Interaction#	Item#	User#	Sparsity
MovieLens	1,000,209	3,706	6,040	95.53%
Pinterest	1,500,809	9,916	55,187	99.73%

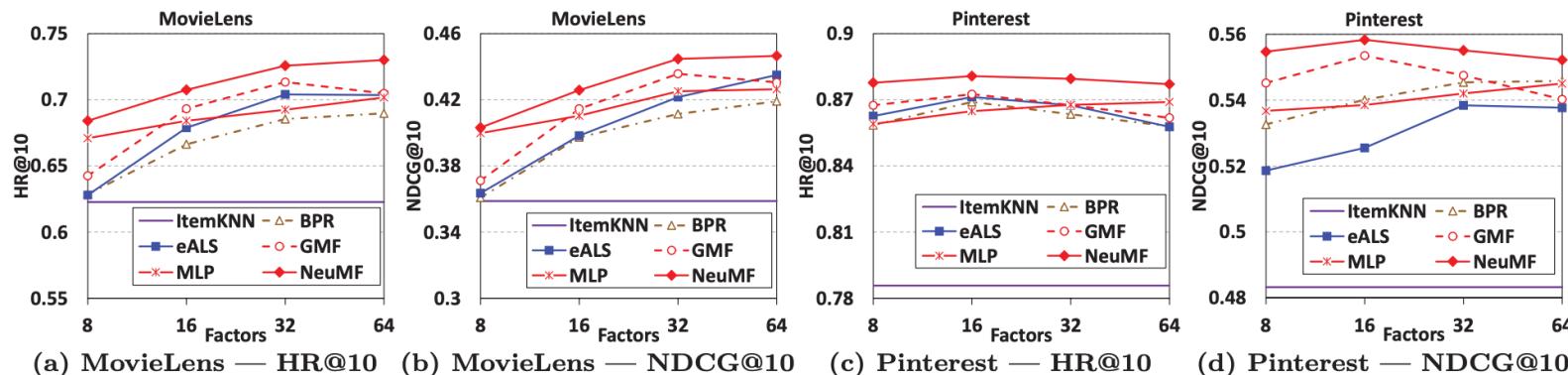


Figure 4: Performance of HR@10 and NDCG@10 w.r.t. the number of predictive factors on the two datasets.

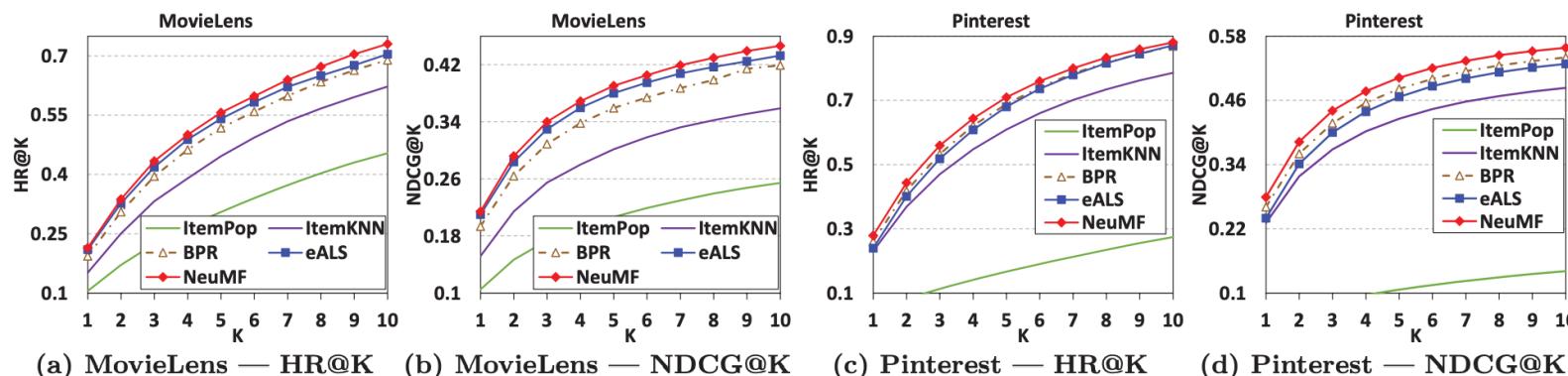


Figure 5: Evaluation of Top-K item recommendation where  $K$  ranges from 1 to 10 on the two datasets.

# NeurMF - Experiments

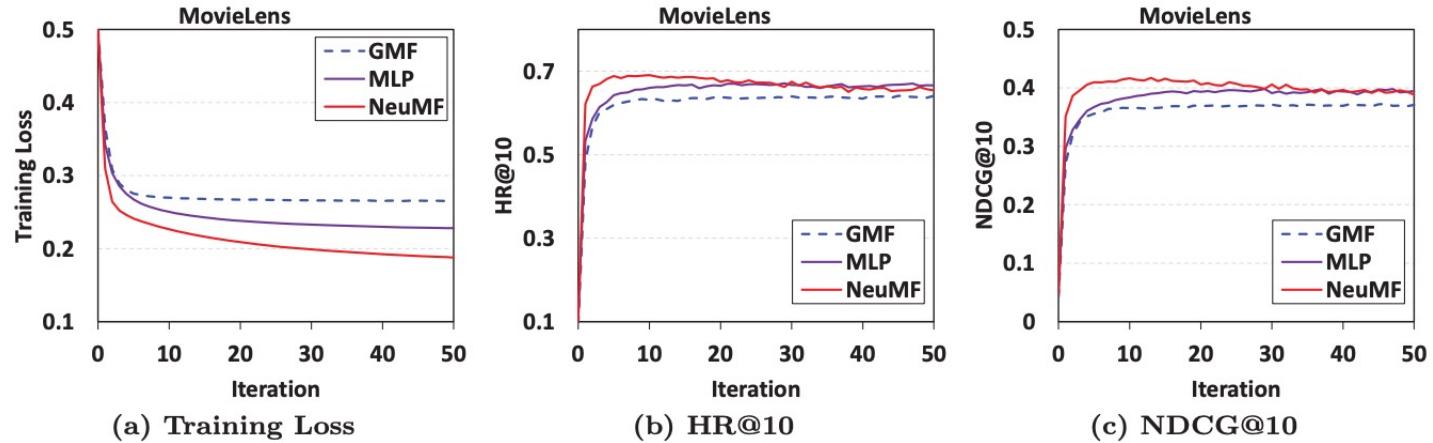


Figure 6: Training loss and recommendation performance of NCF methods *w.r.t.* the number of iterations on MovieLens (factors=8).

**Table 3:** HR@10 of MLP with different layers.

Factors	MLP-0	MLP-1	MLP-2	MLP-3	MLP-4
<b>MovieLens</b>					
<b>8</b>	0.452	0.628	0.655	0.671	<b>0.678</b>
<b>16</b>	0.454	0.663	0.674	0.684	<b>0.690</b>
<b>32</b>	0.453	0.682	0.687	0.692	<b>0.699</b>
<b>64</b>	0.453	0.687	0.696	0.702	<b>0.707</b>
<b>Pinterest</b>					
<b>8</b>	0.275	0.848	0.855	0.859	<b>0.862</b>
<b>16</b>	0.274	0.855	0.861	0.865	<b>0.867</b>
<b>32</b>	0.273	0.861	0.863	<b>0.868</b>	0.867
<b>64</b>	0.274	0.864	0.867	0.869	<b>0.873</b>

## DL helps...

# References

- [1] Learning and Reasoning on Graph for Recommendation, Xiang Wang , Xiangnan He, Tat-Seng Chua , tutorial CIKM 2019
- [2] Yehuda Koren, KDD '09: Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data miningJune 2009 Pages 447–456, <https://doi.org/10.1145/1557019.1557072>
- [3] Pairwise Interaction Tensor Factorization for Personalized Tag Recommendation Steffen Rendle, Lars Schmidt-Thieme, WSDM 2010
- [4] Cao, Jian; Hu, Hengkui; Luo, Tianyan; Wang, Jia; Huang, May; Wang, Karl; Wu, Zhonghai; Zhang, Xing (2015). *Distributed Design and Implementation of SVD++ Algorithm for E-commerce Personalized Recommender System*. Communications in Computer and Information Science. **572**. Springer Singapore. pp. 30–44. [doi:10.1007/978-981-10-0421-6\\_4](https://doi.org/10.1007/978-981-10-0421-6_4). ISBN 978-981-10-0420-9. |
- [5] Deep Matrix Factorization Models for Recommender Systems. HJ Xue, X Dai, J Zhang, S Huang, J Chen - IJCAI, 2017 - ijcai.org
- [6] Sedhain, Suvash, et al. "Autorec: Autoencoders meet collaborative filtering." Proceedings of the 24th international conference on World Wide Web. 2015.
- [7] Rendle, Steffen, et al. "Neural collaborative filtering vs. matrix factorization revisited." Fourteenth ACM Conference on Recommender Systems. 2020.

# Deep Learning for drug design - ARG

- Deep Learning for Drug Discovery
- GNNs for Antibiotic Resistance prediction

# Message Passing Neural Networks for Learning Graph Representations

$$h_1^{(t+1)} = f(W_0^{(t)} h_1^{(t)} + W_1^{(t)} h_2^{(t)} + W_1^{(t)} h_3^{(t)})$$

$$h_2^{(t+1)} = f(W_0^{(t)} h_2^{(t)} + W_1^{(t)} h_1^{(t)} + W_1^{(t)} h_3^{(t)} + W_1^{(t)} h_4^{(t)})$$

$$h_3^{(t+1)} = f(W_0^{(t)} h_3^{(t)} + W_1^{(t)} h_1^{(t)} + W_1^{(t)} h_2^{(t)} + W_1^{(t)} h_4^{(t)})$$

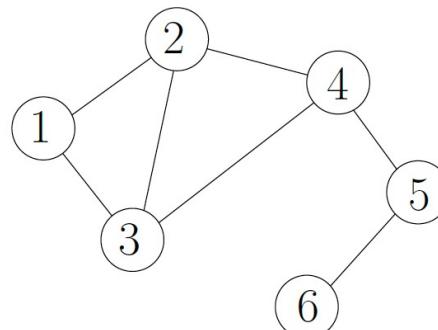
$$h_4^{(t+1)} = f(W_0^{(t)} h_4^{(t)} + W_1^{(t)} h_2^{(t)} + W_1^{(t)} h_3^{(t)} + W_1^{(t)} h_5^{(t)})$$

$$h_5^{(t+1)} = f(W_0^{(t)} h_5^{(t)} + W_1^{(t)} h_4^{(t)} + W_1^{(t)} h_6^{(t)})$$

$$h_6^{(t+1)} = f(W_0^{(t)} h_6^{(t)} + W_1^{(t)} h_5^{(t)})$$

Output of message passing phase:

$$\{h_1^{(T)}, h_2^{(T)}, h_3^{(T)}, h_4^{(T)}, h_5^{(T)}, h_6^{(T)}\}$$



Graph representation:

$$h_G = \frac{1}{6} (h_1^{(T)} + h_2^{(T)} + h_3^{(T)} + h_4^{(T)} + h_5^{(T)} + h_6^{(T)})$$

# Powerful GNN methods

High-dimensional Weisfeiler-Lehman test of isomorphism → a generalization of the WL which colors tuples from  $V^k$  instead of nodes

- $k$ -GNN [Morris et al., AAAI'19]

High-order neighborhoods

- $k$ -hop [Nikolentzos et al., Neural Networks 130]

Approaches that consider the average of all possible permutations of nodes

- RelationalPooling [Murphy et al., ICML'19]

Coloring schemes

- CLIP [Dasoulas et al., IJCAI'20]

Invariant and equivariant linear layers

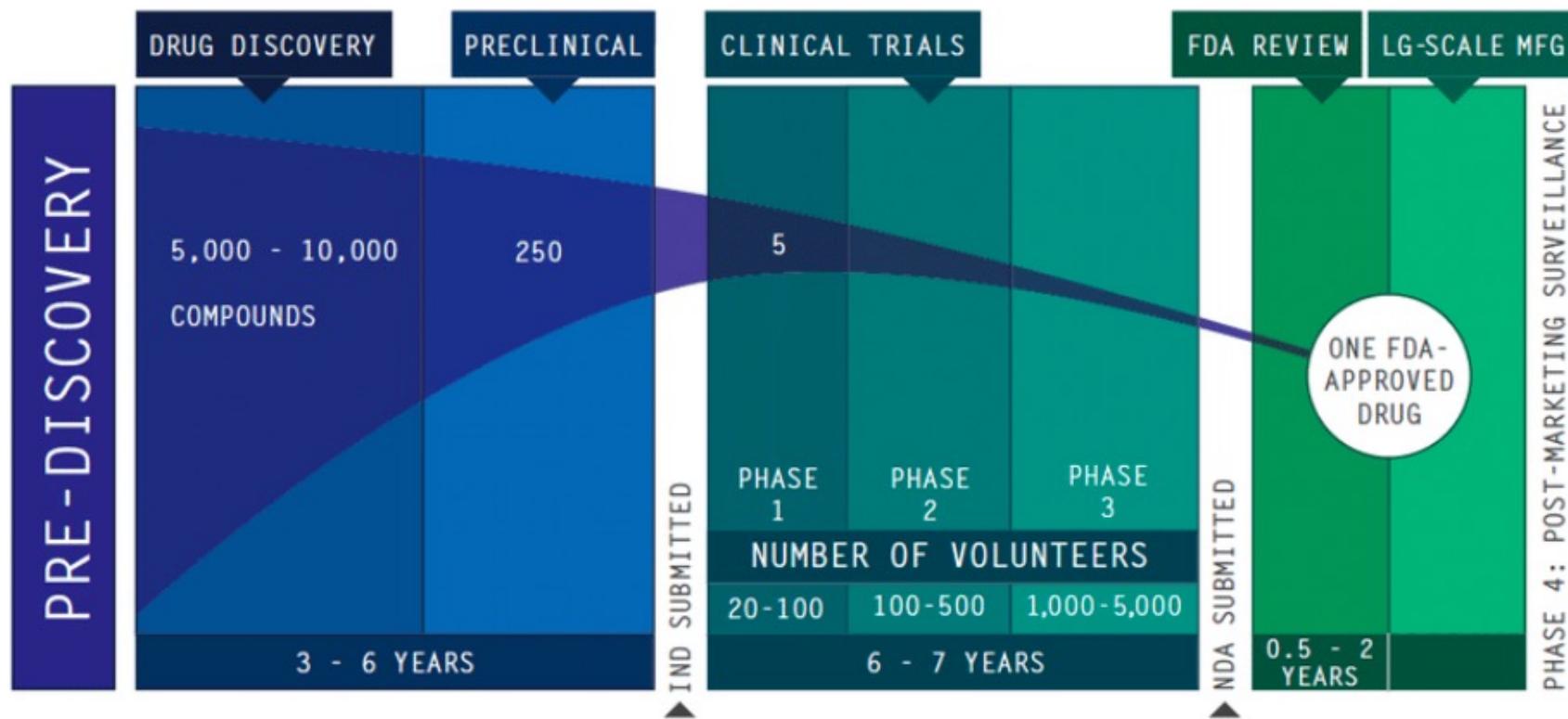
- $k$ -order graph networks [Maron et al., ICLR'19]

# Deep Learning for drug design - ARG

- Deep Learning for Drug Discovery
- GNNs for Antibiotic Resistance prediction

# Drug Discovery – Long process

In average : **10 years** and **\$2.6B**

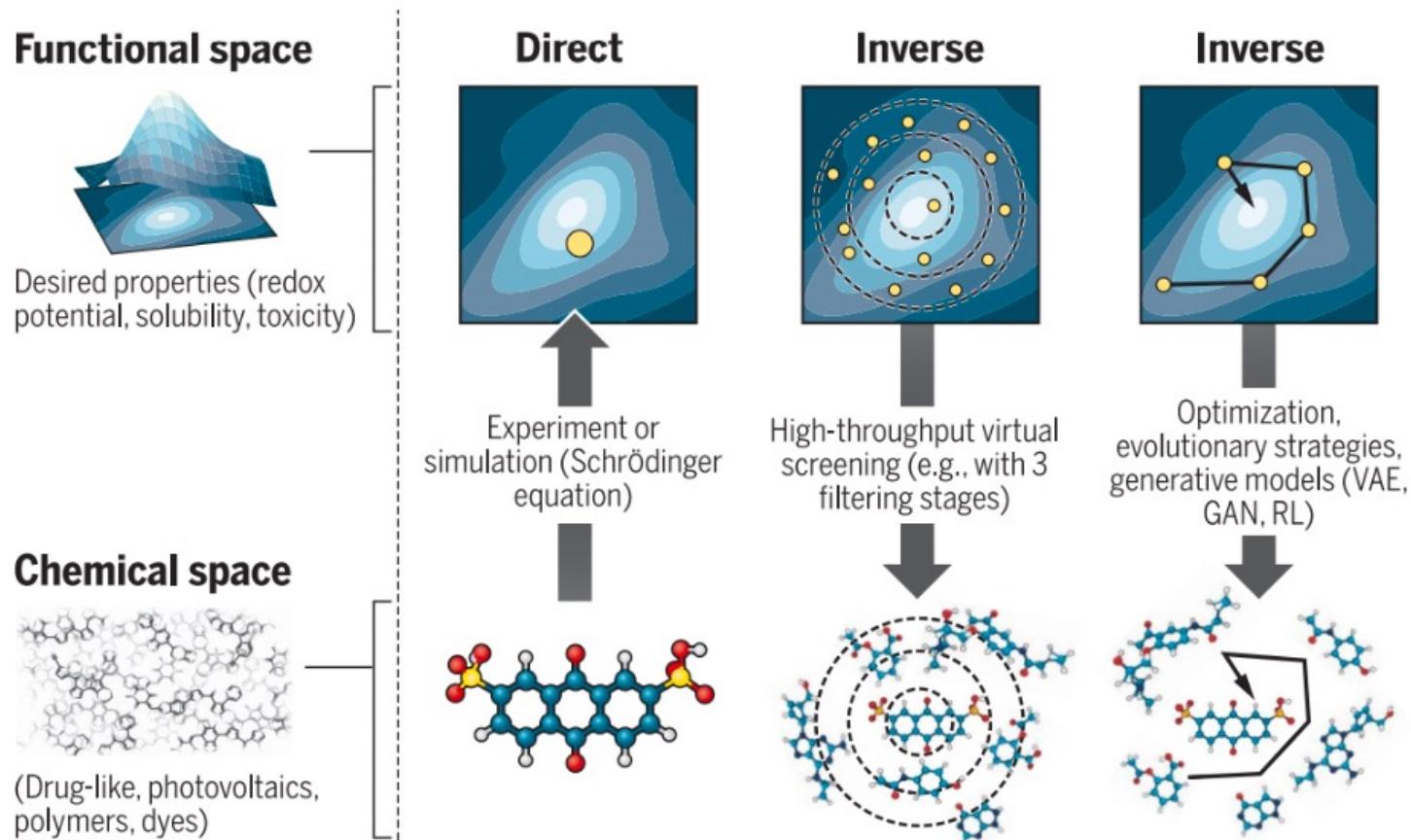


Source : *Pharmaceutical Research and Manufacturers of America*

## Drug discovery is a challenging search problem

- The total number of possible small organic molecules that populate **chemical space** has been estimated to exceed  $10^{60}$ .
- Experimental facilities in industry can only test  $10^5$  compounds/day.  
We need efficient computational methods to design drugs automatically and
- Artificial Intelligence can help !

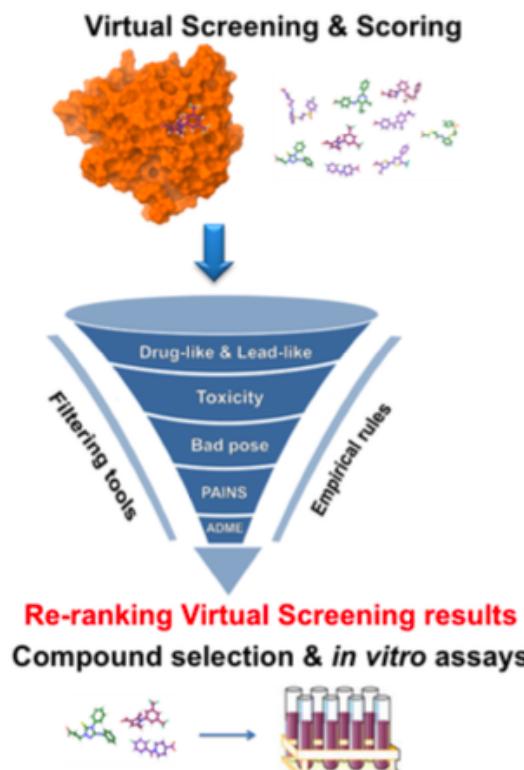
# Approaches towards molecular design



Source : *Sanchez-Lengeling et al., Science 361, 360–365 (2018)*

# Virtual Screening

- **Virtual screening** : assess whether a compound is a good drug using computation models (Walters et al., 1998 ; McGregor et al., 2007 ; ...)

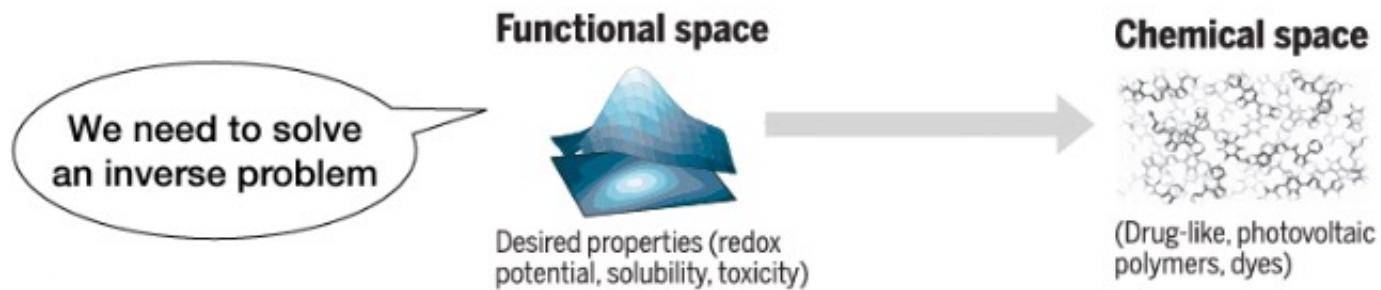
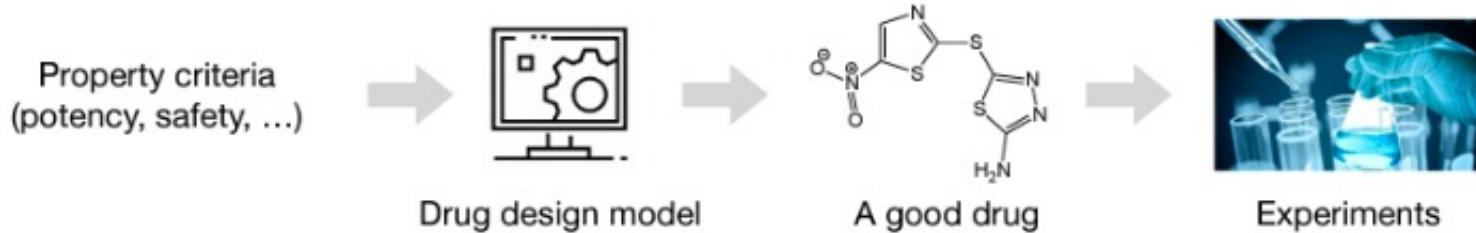


- Virtual screening is much faster than experimental screening in web labs.
- It can test  $10^8$  compounds within a day, while experimental screening takes years.
- It is also much cheaper than experimental screening.

Source : *Surabhi, Surabhi Singh, BK. (2018).*

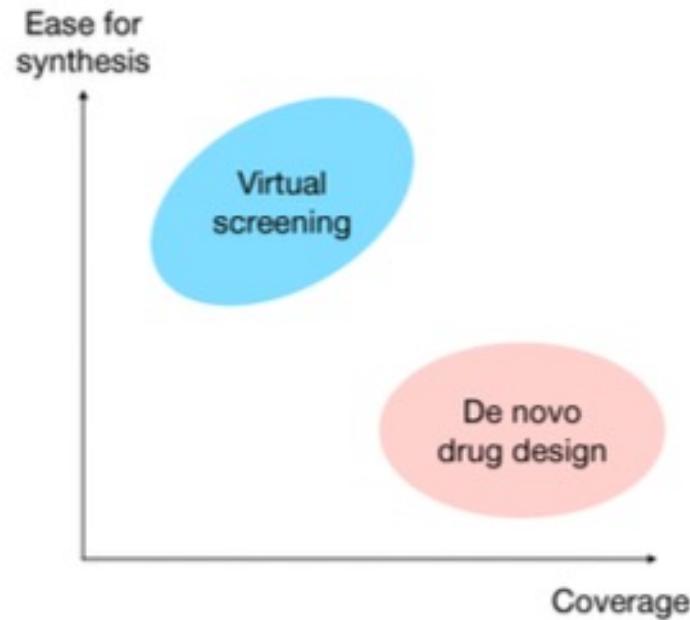
# De Novo Drug Design

- **De novo drug design** : directly generate a compound with desired properties  
(Moon et al., 1991 ; Clark et al., 1995 ; Schneider Fechner, 2005 ; ...)



Source : Wengong Jin, Deep Learning for drug discovery (2021)

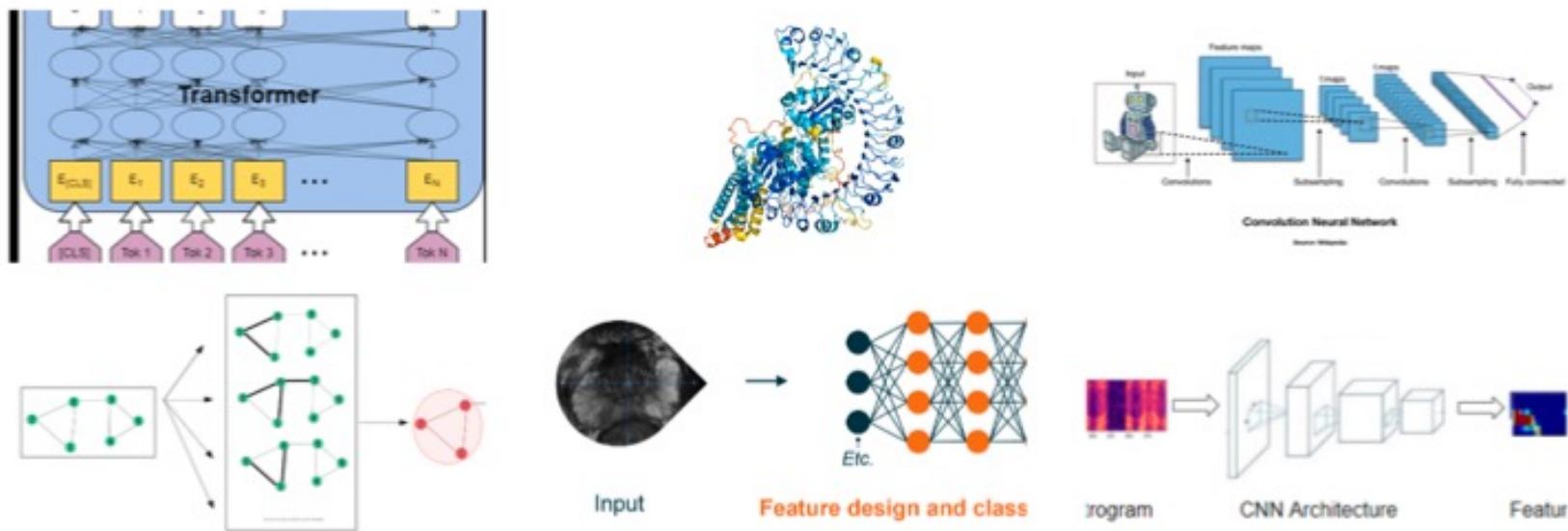
# Virtual Screening vs De Novo Drug Design



- Virtual screening is restricted to commercially available compounds (e.g., ZINC library - containing  $2 * 10^8$  compounds).
- De novo drug design can explore the entire chemical space efficiently.
- Issue : feasibility of the new compounds synthesis.
- Traditional Techniques for both based on hand-crafted features and hand-designed rules
- Source : *Wengong Jin, Deep Learning in drug discovery (2021)*

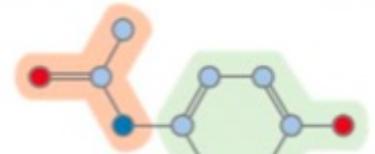
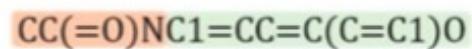
# Deep Learning – a promising direction

- Deep Learning has achieved many advances in recent years, with human-level accuracy in many tasks (Image Processing, Natural Language Processing, Recommendation Systems, Protein Folding, ...)



# Molecules representations

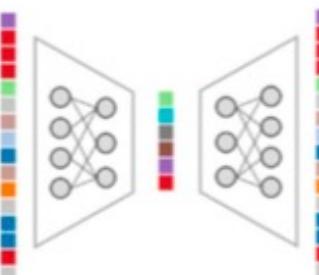
1) SMILES



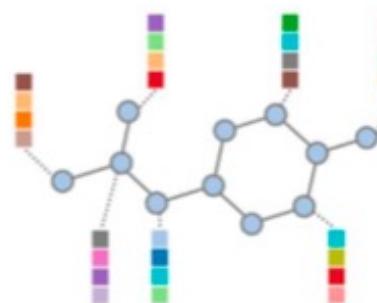
2) Fingerprint



3) Learned feature from AE



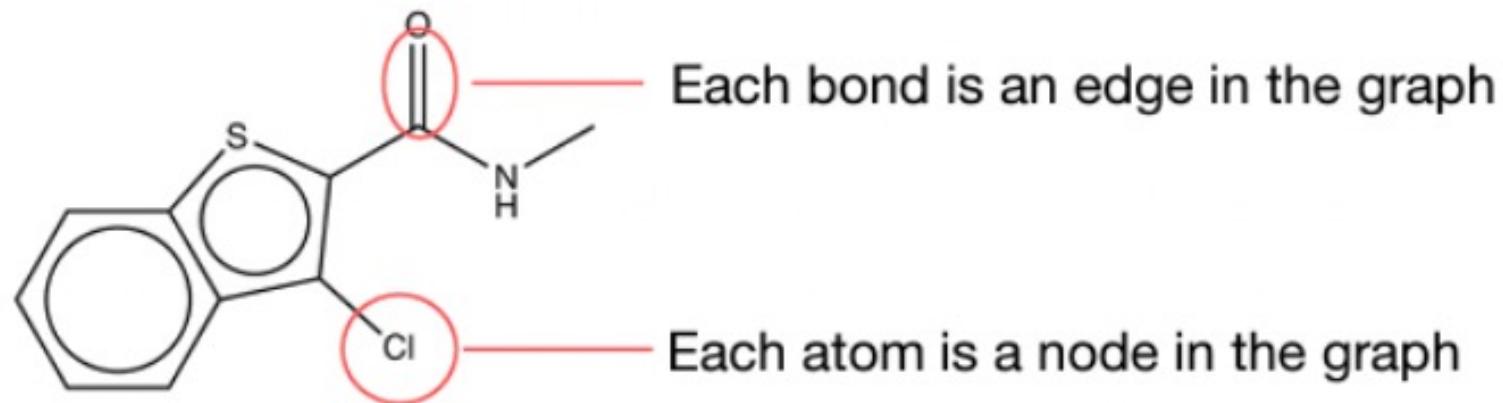
5) Molecular graph



Source : *Kim, Jintae et al.(2021)*

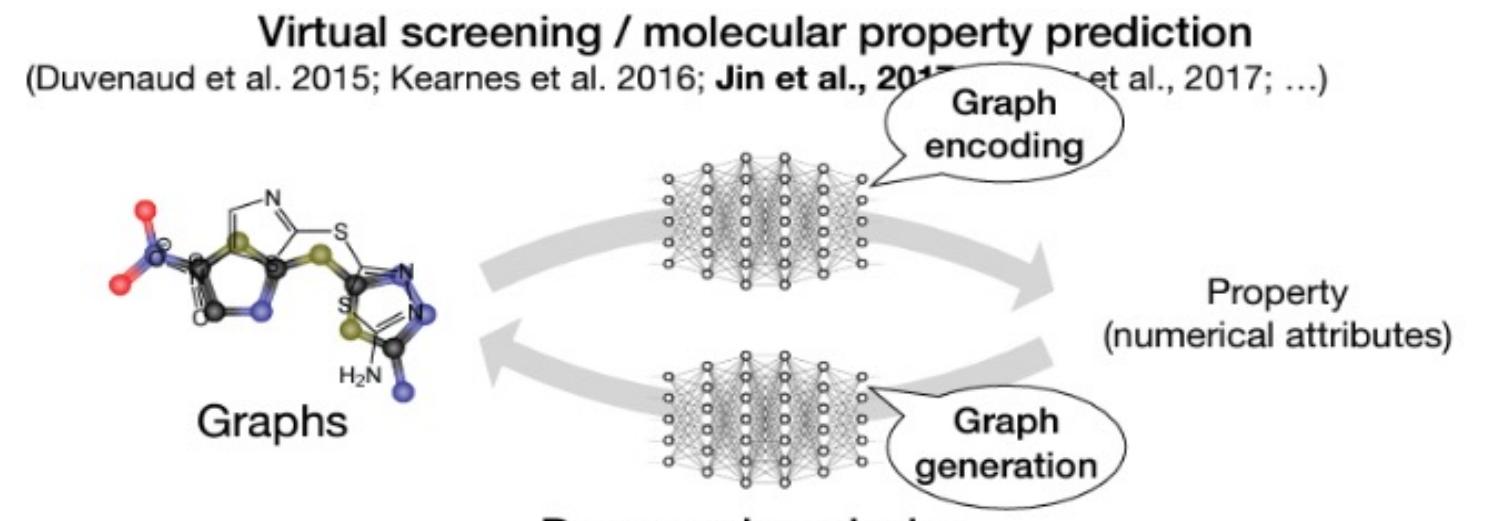
# Molecules as Graphs

- Molecules can be represented as graphs.
- Node features : Atom type, Formal charge, etc.
- Edge features : Bond type, etc.



# Deep Learning for drug discovery

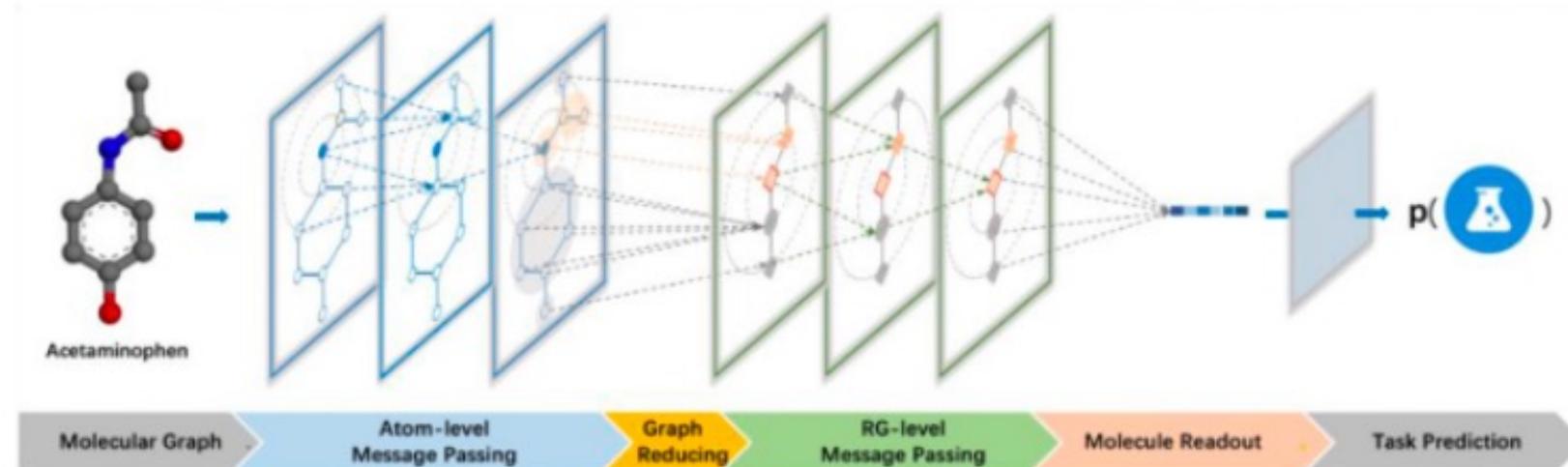
- **Virtual Screening** : Can be seen as a classification problem of molecules.
- **De Novo Drug Design** : Can be seen as a generative problem of molecules, with constraints.



Source : *Wengong Jin, Deep Learning for drug discovery (2021)*

# Graph Neural Networks for molecules

- Rich history of GNNs (Gori et al., 2005, Scarselli et al., 2009, Duvenaud et al. 2015, Kearnes et al. 2016, Jin et al., 2017, Gilmer et al., 2017, Zitnik et al., 2018, etc.)
- Adapted to leverage the structure and properties of graphs.
- Used in Chemistry, Natural Language Processing, Recommendation systems, Computer Vision, Graph Generation ,etc.



Source : Kong, Y., Zhao, X., Liu, R. et al., 2022

# Deep Learning for drug design - ARG

- Deep Learning for Drug Discovery
- GNNs for Antibiotic Resistance prediction

# Structure-Aware Antibiotic Resistance Classification Using Graph Neural Networks [1]

**Definition** : Antimicrobial resistance (AR) happens when germs like bacteria and fungi develop the ability to defeat the drugs designed to kill them.

## Why ?

- Major Health Challenge
- Natural Process, but enhanced by overuse and misuse in humans and animals
- Increased mortality (up to 10 million people by 2050)[5]
- Few antibiotics are developed each year, due to expensive process and low profit.

**Antibiotics** : Many antibiotics work by inactivating an essential bacterial protein.

**How** ? It can happen for many reasons :

- Spontaneous mutation in the bacterium's DNA, which change the essential bacterial protein.
- Transfer of antibiotic-resistance genes

[1] A. Qabel et. al, "Structure-Aware Antibiotic Resistance Classification Using Graph Neural Networks", <https://www.biorxiv.org/content/10.1101/2022.10.06.511103v1>, accepted in the AI4Science WS@Neurips2022

## ARG prediction - Current baselines

Sequence based Algorithms :

- BLAST: Finds regions of local similarity between sequences.
- DIAMOND : Another sequence alignment tool, faster than BLAST.

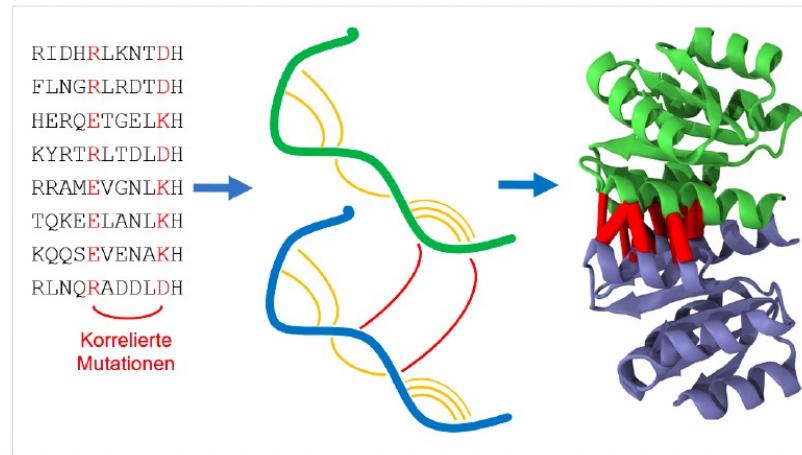
Machine Learning :

- DeepARG: Features = BLAST best hit score with some known ARs.
- TRAC: Pretraining of an LSTM-based LM + finetune on AR prediction task.
- ARG-SHINE: an ensemble method using BLAST results and deep convolutional network on the sequences.

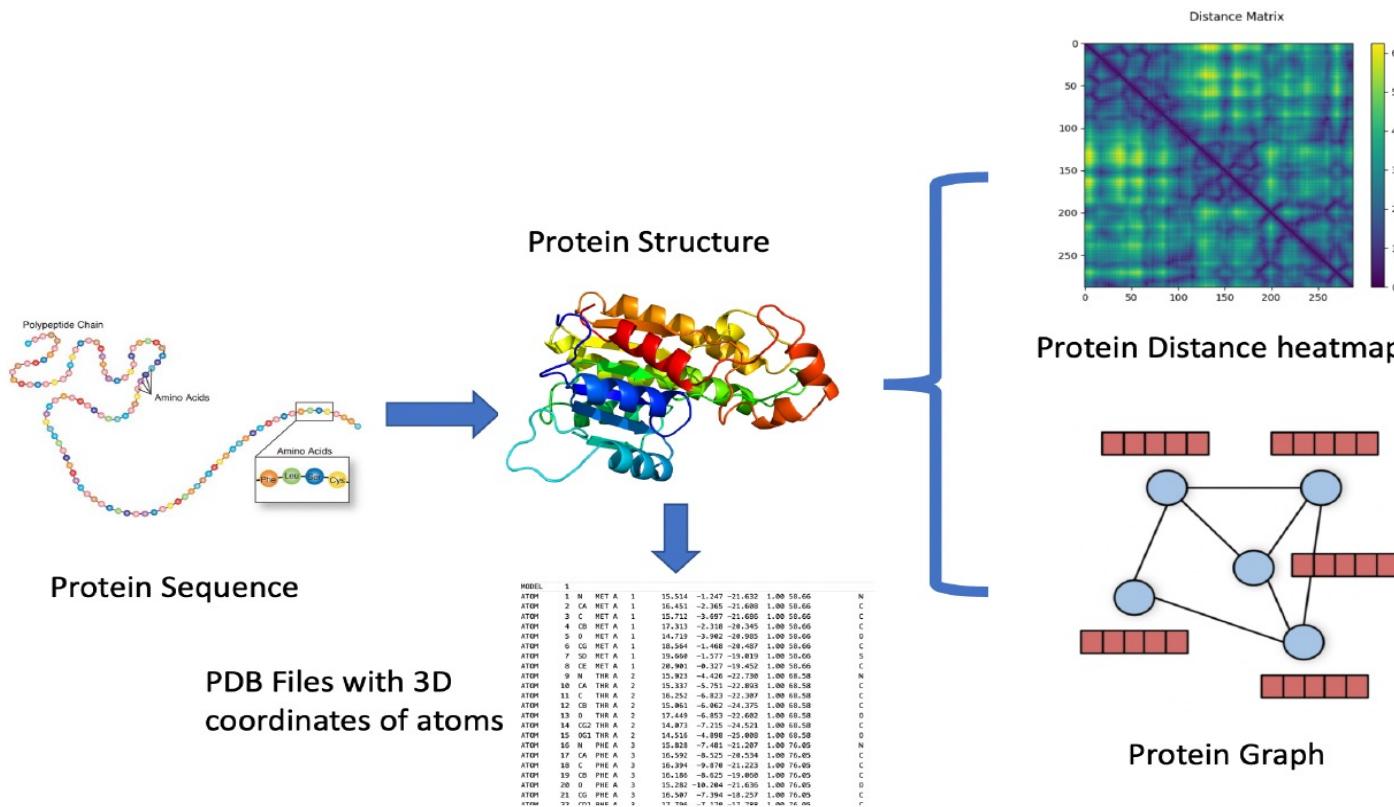
All these methods use Protein Sequence !

# Protein folding

- Protein Folding: physical process - protein chain translated to its native three-dimensional structure – *a very tough problem*
- computational methods using AI to predict the protein structure:
  - DeepMind developed *AlphaFold2*, a model that achieved high quality *structure prediction*.
  - Facebook released *ESM2*, a Large Language model pretrained on sequences that achieves similar results to AlphaFold when using the structure module of AlphaFold2.



# How to use protein structure

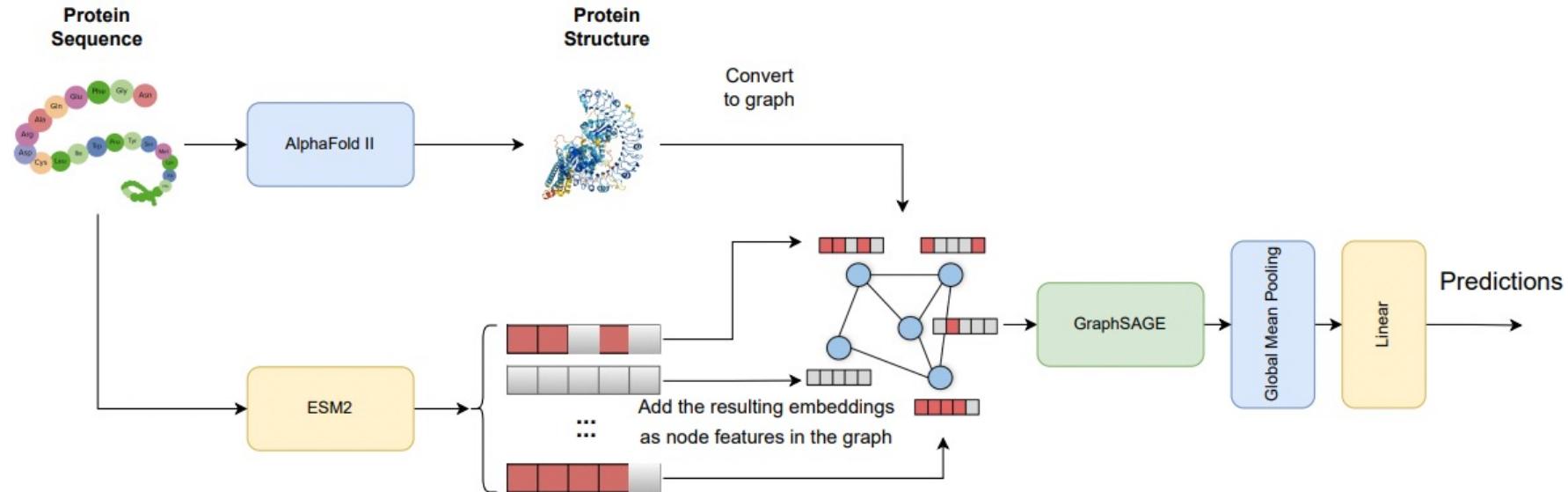


## Protein 3D structure - Features

Proteins represented as Graphs, where :

- **Nodes** : Amino-Acids
- **Node Features** :
  - ▶ **Amino-Acid type** : One-hot-encoding of the type (20 different types).
  - ▶ **Phi** angle : phi-angles of each residue.
  - ▶ **Psi** angle : psi-angles of each residue.
  - ▶ **RSA** : relative solvent accessibility of each residue.
- **Edges** : defined based on distance threshold ( $<12 \text{ \AA}$ ) + Peptide bonds.
- **Edge Features** :
  - ▶ **Distance** : distance between residues.
  - ▶ **Bond type** : Peptide bond or distance bond

# LM + GNN for AR prediction



The framework consists of two main components : (1) The **AlphaFold** model predicting the 3D structure of an input protein; (2) A classifier, consisting of a GraphSAGE model applied to the graph constructed from the structure and with node features from a pretrained protein language model (**ESM2**) applied to the sequence.

# ARG Data set

- **COALA** : COllection of ALI Antibiotic resistance gene databases
- 70% max identity between sequence pairs.
- 16 AR classes.
- **Splitting :**
  - ▶ train set : 80% of the dataset
  - ▶ test set : 20% of the dataset
  - ▶ Stratified Splitting between train and test set.
  - ▶ Ensure test sequences have at most 40% identity with the training set sequences : To do so, we use CD-HIT, which produces global identity score between sequence pairs and cluster them and iteratively add clusters, making sure we have a stratified split.
  - ▶ → 5 different splits of data

# GNN for AR prediction

Results on test set				
	Accuracy	Accuracy (<50% id)	Accuracy (>50% id)	Accuracy (Homologue not found)
BLAST	0.6542±0.0057	0.7511±0.0178	<b>0.9529±0.0118</b>	0
Diamond	0.5886±0.0062	0.6469±0.0172	<b>0.9511±0.0153</b>	0.0005±0.0010
TF-IDF	0.5419±0.0162	0.5518% (± 0.172%)	0.8405% (± 0.207%)	0.3543% (± 0.125%)
TRAC	0.6980±0.0066	0.7278±0.0177	0.9092±0.0216	0.3683±0.0436
ARG-SHINE-CNN	0.6834±0.0127	0.6960±0.0178	0.9288±0.0070	0.3800±0.0388
<b>Efficient-Net</b>	0.6983±0.0089	0.7310±0.0217	0.9167±0.0236	0.3507±0.0288
<b>GraphSAGE</b>	0.6666±0.0098	0.6938±0.0168	0.8972±0.0148	0.3257±0.0318
<b>ESM2_8M</b>	0.7086±0.0133	0.7437±0.0170	0.9260±0.0050	0.3535±0.0493
<b>ESM2_650M</b>	<b>0.7310±0.0138</b>	<b>0.7652±0.0174</b>	0.9094±0.0080	<b>0.4216±0.0357</b>
<b>ESM2_8M_GraphSAGE</b>	<b>0.7290±0.0065</b>	<b>0.7649±0.0130</b>	0.9300±0.011	<b>0.3890±0.0398</b>

## Conclusions

- Overall improvement of accuracy : More than 3%.
- Transformer Language models perform very well.
- Combining Sequence + Structure information leads to better performances.
- Deep Learning methods showed good performance detecting Ars even without similar homologues in the training set.

## Vision – structure assisted protein pretrained models

- NLP pretrained models
  - GPT3, BARTHEZ
  - Very successful for downstream and generative tasks (i.e. translation, summarization, text generation...)
- Pretrain protein models based on structure as well
  - Capitalise on Alphafold2 database (~200M proteins) and large LM for proteins models (ESM2)
  - Downstream tasks (i.e. protein function prediction)
  - Generative tasks (new proteins ...)
- multimodal pretraining of protein structures and textual descriptions (functions)

# THANK YOU!

- DASCIM web page:  
<http://www.lix.polytechnique.fr/dascim>
- Google Scholar: <https://bit.ly/2rwmvQU>
  - Twitter: @mvazirg