

# Workshop #2

Amaru Cuba Gyllensten  
[amaru.cuba.gyllensten@ri.se](mailto:amaru.cuba.gyllensten@ri.se)

Ariel Ekgren  
[ariel.ekgren@ri.se](mailto:ariel.ekgren@ri.se)

# Översikt

1. Introduktion
2. Presentation
3. NLU-Colab
4. Pilot smartflow

# Språkmodeller

Amaru Cuba Gyllensten

[amaru.cuba.gyllensten@ri.se](mailto:amaru.cuba.gyllensten@ri.se)

Ariel Ekgren

[ariel.ekgren@ri.se](mailto:ariel.ekgren@ri.se)

# Översikt

1. Vad är NLU?
2. Hur görs NLU?
  - a. Transformers
  - b. Förträning
  - c. Finetuning
3. NLU i praktiken

1. Vad är NLU?
2. Hur görs NLU?
  - a. Transformers
  - b. Förträning
  - c. Finetuning
3. NLU i praktiken

# Natural Language Understanding

NLU är en gren inom ML för att lösa en mängd språkförståelse-relaterade problem:

- Klassifikation
- Sammanfattning
- Namnigenkänning
- Översättning

# Natural Language Understanding

Idag används förtränade språkmodeller till nästan alla NLU-problem.

Vad är en språkmodell?

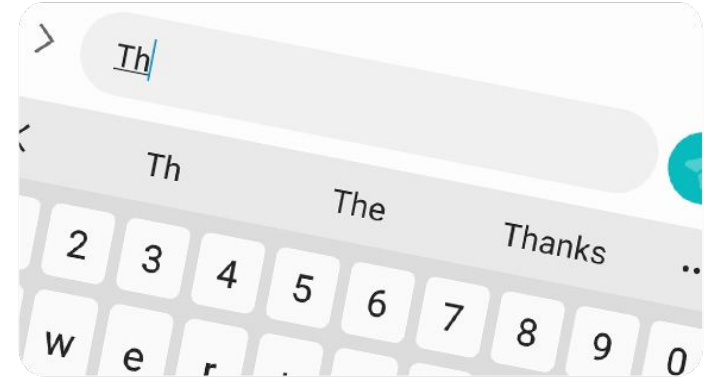
# Språkmodeller i korthet

Abstrakt: Språkmodeller är modeller som lär sig språkets *distribution*.

Konkret: Språkmodeller är modeller som “gissar” ord, givet en kontext.



# Språkmodeller i korthet



- Gissa nästa ord:

“.. ett höganäskrus i en svångrem om ???”

- Gissa ord i kontext:

“... ett ??? i en svångrem om halsen ...”

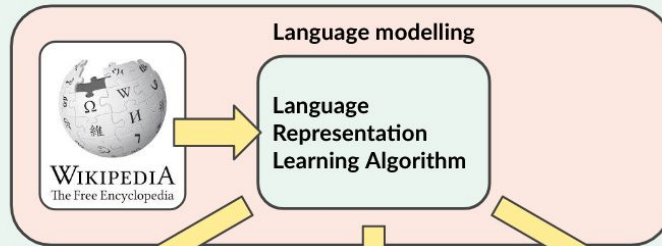
# Språkmodeller i korthet

Språkmodellen i sig är sällan intressant.

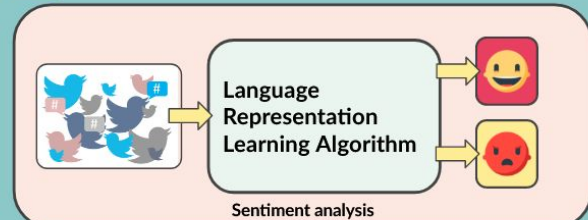
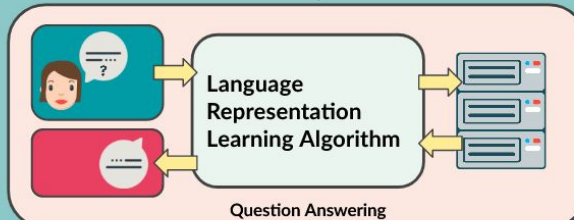
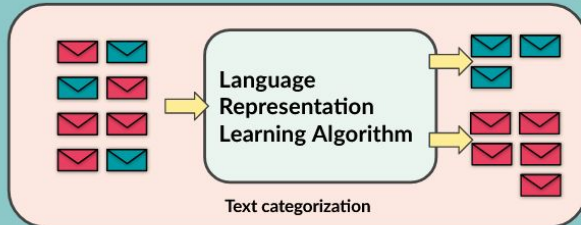
Språkmodellen är grunden med vilket vi bygger NLU-lösningar!

# Varför språkmodeller?

## Phase 1 - Unsupervised pre-training



## Phase 2 - Fine tune to tasks



# Varför språkmodeller?

Språkförståelse är *nödvändigt* för dokument-klassificering, fråga-svarssystem, sentimentanalys et.c.

Språkförståelse är *nödvändigt* för språkmodellering.

En tillräckligt bra språkmodell har lärt sig “förstå språk” till den grad att den enklare kan anpassas till mer specifika problem.

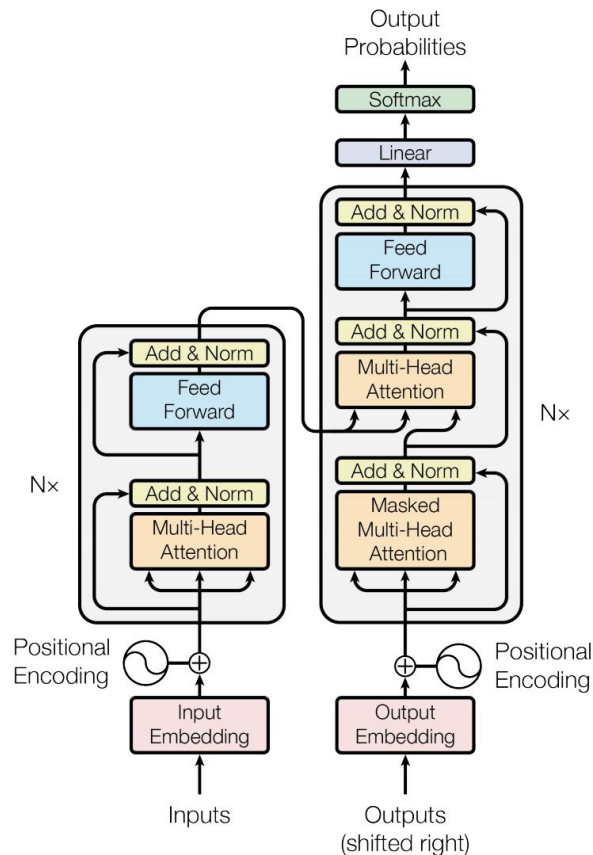
1. Vad är NLU?
2. Hur görs NLU?
  - a. Transformers
  - b. Förträning
  - c. Finetuning
3. NLU i praktiken

# NLU ❤️ Transformers

Den “bästa” familjen av språkmodeller för tillfället är Transformers.

Transformer-modellen introducerades 2019 med artikeln Attention is All You Need.

Vaswani, Ashish, et al. "Attention is all you need." Advances in neural information processing systems. 2017.



# Text -> Tokens

För att hantera text behöver man konvertera den till något ett neuralt nätverk kan använda.

Detta görs genom tokenisering. Tokenisering styckar upp en text i delar, eller tokens. t.ex. “en katt satt på en hatt” -> [en, katt, satt, på, en, hatt].

Tokenisering ger oss också ett vokabulär av alla tokens som förekommer i datan.

# Tokens -> Embeddings

Efter tokenisering har vi ett vokabulär av tokens.

För varje token skapar vi en “embedding”. Detta är en parameteriserad vektor, som stoppas in som input i neurala nätverk.

Embedding Matrix

“the”:						
“be”:						
“to”:						
“of”:						
“and”:						
“a”:						
“in”:						
“that”:						

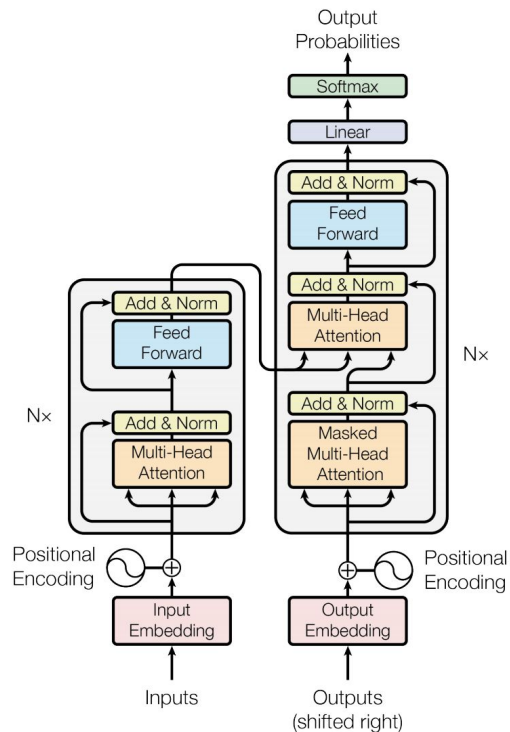
vocabulary size x embedding dimension



# Embeddings -> Representation

Embeddings stoppas sedan in som input i själva transformern.

Output är omvandlade embeddings.



"En katt satt på en hatt."								
<B>	En	katt	satt	på	en	hatt	.	<E>
2	59	7808	1549	68	59	9226	7	3
↗	↘	↖	↓	↙	↗	↖	↖	↗
→	↘	↗	↓	↙	→	↖	↓	↗
⋮								
↗	↘	←	↓	↙	→	↗	→	↗

1. Vad är NLU?
2. Hur görs NLU?
  - a. Transformers
  - b. Förträning**
  - c. Finetuning
3. NLU i praktiken

# Förträning

Förträning går ut på att en låter algoritmen träna på en stor mängd icke annoterad textdata.

Till exempel text från Kungliga Biblioteket och wikipedia.

Detta görs genom att gissa ord i kontext.



Källa: <https://xkcd.com/1838/>

"En katt satt på en hatt."								
<B>	En	katt	satt	på	en	hatt	.	<E>
2	59	7808	1549	68	59	9226	7	3
↗	↘	↖	↓	↙	↗	↖	↖	↗
→	↘	↗	↓	↙	→	↖	↓	↗
⋮								
↗	↘	←	↓	↙	→	↗	→	↗

"En katt satt på en hatt."								
<B>	En	???	satt	???	en	hatt	.	<E>
2	59	1	1549	1	59	9226	7	3
↗	↘	↑	↘	↑	↘	↖	↖	↑
→	↘	↗	↘	↗	→	↖	↓	↖
⋮								
↗	↘	←	↘	↗	↘	↖	→	↗

"En katt satt på en hatt."								
<B>	En	???	satt	???	en	hatt	.	<E>
2	59	1	1549	1	59	9226	7	3
↗	↘	↑	↘	↑	↘	↖	↖	↑
→	↘	↗	↘	↗	→	↖	↓	↖
⋮								
↗	↘	←	↘	↗	↘	↖	→	↗



Bra!



Dåligt!

1. Vad är NLU?
2. **Hur görs NLU?**
  - a. Transformers
  - b. Förträning
  - c. Finetuning**
3. NLU i praktiken



# Fine-tuning

När man har en förtränad modell är det vanligaste att man finjusterar modellerna mot ett annoterat dataset. Med det förfarandet fortsätter man uppdatera modellerna mot ett specifikt mål.

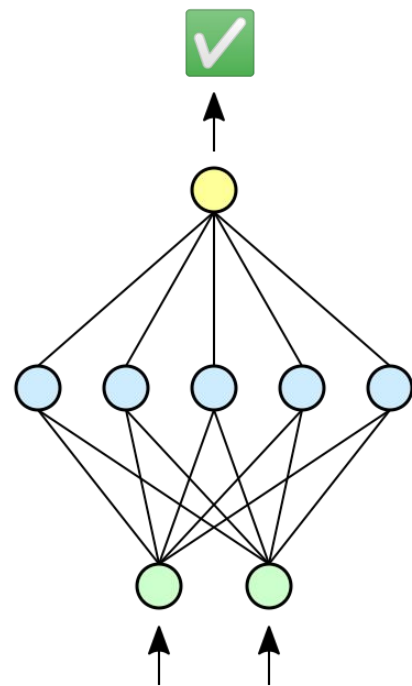
Det kräver dock tusentals annoterade exempel för att uppnå bra resultat.

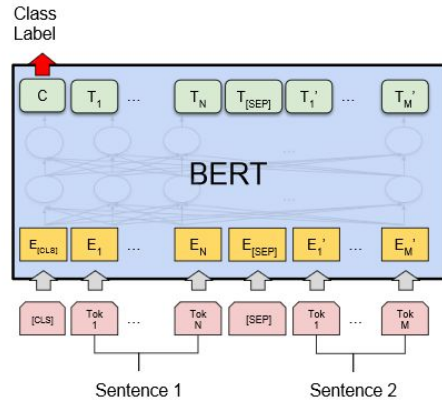
## Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.

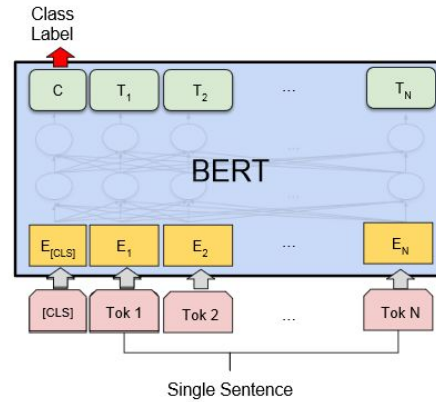


Korrekt : "En katt satt på en hatt."								
<B>	En	katt	satt	på	en	hatt	.	<E>
2	59	7808	1549	68	59	9226	7	3
↗	↘	↖	↙	↘	↗	↖	↙	↗
→	↘	↗	↓	↘	→	↖	↓	↗
⋮								
↗	↘	↖	↙	↘	↗	↖	↙	↗

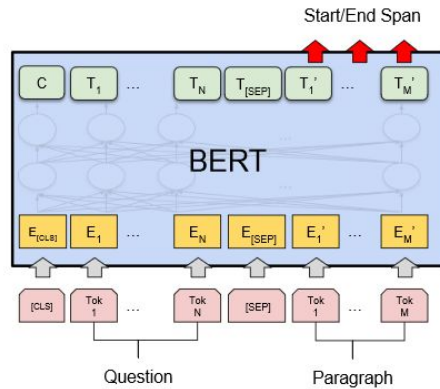




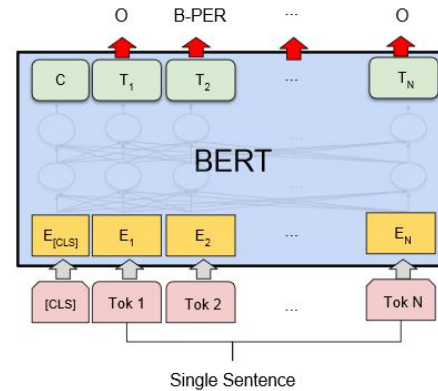
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA



(c) Question Answering Tasks:  
SQuAD v1.1

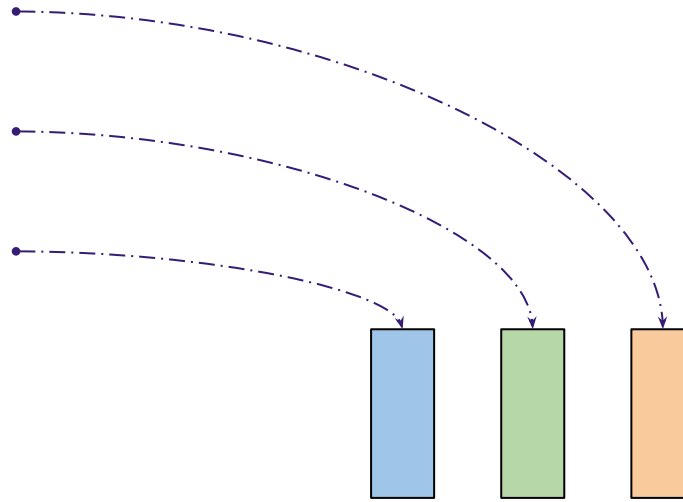


(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

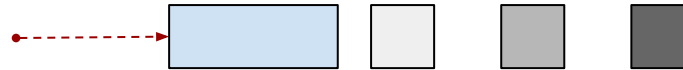
Hej! Återställ ditt lösenord via ...

Hej! Om du har ett öppet dokument ...

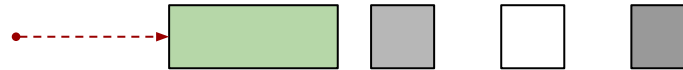
Hej! Gå till receptionen och anmäl ...



Hur återställer jag mitt lösenord?



Hur sparar jag filer i teams?



Jag har tappat bort min nyckel!



1. Vad är NLU?
2. Hur görs NLU?
  - a. Transformers
  - b. Förträning
  - c. Finetuning
3. **NLU i praktiken**

# Teckenkodning

Det finns många textkodningar:

ascii, latin-1, windows-1250, **unicode** (utf-8)

Unicode är bra (och standard)! men: **Normalisera**.

Source		NFD		NFC	Source		NFD		NFC
Å 212B	:	Å ◌ 0041 030A		Å 00C5	Å 00C5	:	Å ◌ 0041 030A		Å 00C5
Ω 2126	:	Ω 03A9		Ω 03A9	Ô 00F4	:	Ô ◌ 006F 0302		Ô 00F4

	_0	_1	_2	_3	_4	_5	_6	_7	UTF-8
(1 byte) 0_	NUL 0000	SOH 0001	STX 0002	ETX 0003	EOT 0004	ENQ 0005	ACK 0006	BEL 0007	
(1) 1_	DLE 0010	DC1 0011	DC2 0012	DC3 0013	DC4 0014	NAK 0015	SYN 0016	ETB 0017	
(1) 2_	SP 0020	! 0021	" 0022	# 0023	\$ 0024	% 0025	& 0026	' 0027	
(1) 3_	0 0030	1 0031	2 0032	3 0033	4 0034	5 0035	6 0036	7 0037	
(1) 4_	@ 0040	A 0041	B 0042	C 0043	D 0044	E 0045	F 0046	G 0047	
(1) 5_	P 0050	Q 0051	R 0052	S 0053	T 0054	U 0055	V 0056	W 0057	
(1) 6_	` 0060	a 0061	b 0062	c 0063	d 0064	e 0065	f 0066	g 0067	
(1) 7_	p 0070	q 0071	r 0072	s 0073	t 0074	u 0075	v 0076	w 0077	
8_	• +00	• +01	• +02	• +03	• +04	• +05	• +06	• +07	
9_	• +10	• +11	• +12	• +13	• +14	• +15	• +16	• +17	
A_	• +20	• +21	• +22	• +23	• +24	• +25	• +26	• +27	
B_	• +30	• +31	• +32	• +33	• +34	• +35	• +36	• +37	
(2) C_	2 0000	2 0040	LATIN 0080	LATIN 00C0	LATIN 0100	LATIN 0140	LATIN 0180	LATIN 01C0	
(2) D_	CYRIL 0400	CYRIL 0440	CYRIL 0480	CYRIL 04C0	CYRIL... 0500	ARMENI 0540	HEBREW 0580	HEBREW 05C0	
(3) E_	INDIC 0800	MISC. 1000	SYMBOL 2000	KANA... 3000	CJK 4000	CJK 5000	CJK 6000	CJK 7000	
(4) F_	SMP... 10000	SPUA-B 40000	SSP... 80000	SSP... C0000	SPUA-B 100000	4 140000	4 180000	4 1C0000	

# Textdataformat

Text kan finnas i många format! PDF, Word, xml, et.c.

PDF: Jättebra utdataformat! Väldigt dåligt indataformat!

Som indata vill man helst ha strukturerad text i **unicode**. T.ex. json, csv, xml.

# Huggingface



# Transformers

build **passing** license **Apache-2.0** website **online** release **v2.0.0**

De facto standard för transformerlösningar. Huggingface är en open source lösning för språkmodeller av typen Transformers. Fungerar med både tensorflow och pytorch. De erbjuder också lösningar för att ladda ner förtränade modeller, tokenisering samt dataset.

Om man vill applicera språkmodeller snabbt och smidigt är huggingface ett bra alternativ. <https://huggingface.co/>



# Vilken språkmodell ska man välja?

Välj en modell någon annan tränat. Större och bättre språkmodeller kommer komma.

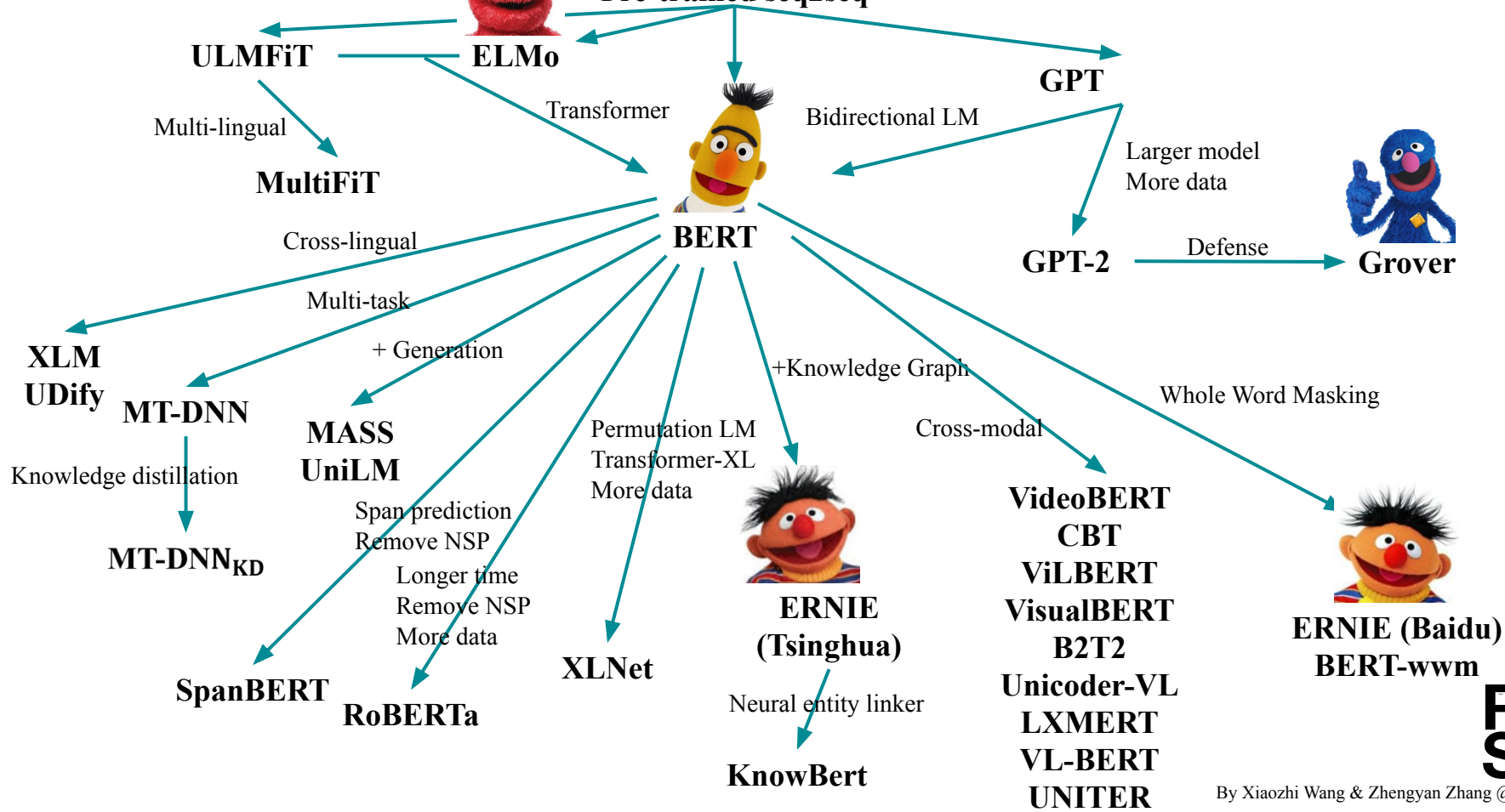
Det viktiga är att det finns finetuning data!

- Finns textmängder tillgängliga som kan användas?
- Finns det processer som skapar annoterad, eller kopplad text?

# Semi-supervised Sequence Learning

context2Vec

Pre-trained seq2seq



# Kurser och läromaterial

Allmänna pytorch tutorials

<https://pytorch.org/tutorials/>

En numera klassisk online-kurs i maskininlärning

<https://www.coursera.org/learn/machine-learning>

En kurs i NLP med huggingface biblioteken (mycket användbar)

<https://huggingface.co/course/chapter1>

# Djupdykning i Transformers

<https://arxiv.org/abs/1706.03762>

<http://jalammar.github.io/illustrated-transformer/>

<https://github.com/huggingface/transformers>

# Verktvg



<https://git-scm.com/>



<https://www.docker.com/>



<https://docs.conda.io/>





<https://pytorch.org/>



**HUGGING FACE**

<https://huggingface.co/>



# Weights & Biases

<https://wandb.ai/>

# vi - terminal text editor

```
ddos@DESKTOP-UVQDIBV: ~  
ddos@DESKTOP-UVQDIBV:~$ vi --help  
VIM - Vi IMproved 8.0 (2016 Sep 12, compiled Jun 06 2019 17:31:41)  
  
usage: vim [arguments] [file ..]      edit specified file(s)  
or: vim [arguments] -                read text from stdin  
or: vim [arguments] -t tag            edit file where tag is defined  
or: vim [arguments] -q [errorfile]    edit file with first error  
  
Arguments:  
--                Only file names after this  
-v                Vi mode (like "vi")  
-e                Ex mode (like "ex")  
-E                Improved Ex mode  
-s                Silent (batch) mode (only for "ex")  
-d                Diff mode (like "vimdif")  
-y                Easy mode (like "evim", modeless)  
-R                Readonly mode (like "view")  
-Z                Restricted mode (like "rvim")  
-m                Modifications (writing files) not allowed  
-M                Modifications in text not allowed  
-b                Binary mode  
-l                Lisp mode  
-C                Compatible with Vi: 'compatible'  
-N                Not fully Vi compatible: 'nocompatible'  
-V[N][fname]      Be verbose [level N] [log messages to fname]  
-D                Debugging mode  
-n                No swap file, use memory only  
-r                List swap files and exit  
-r (with file name) Recover crashed session  
-L                Same as -r
```