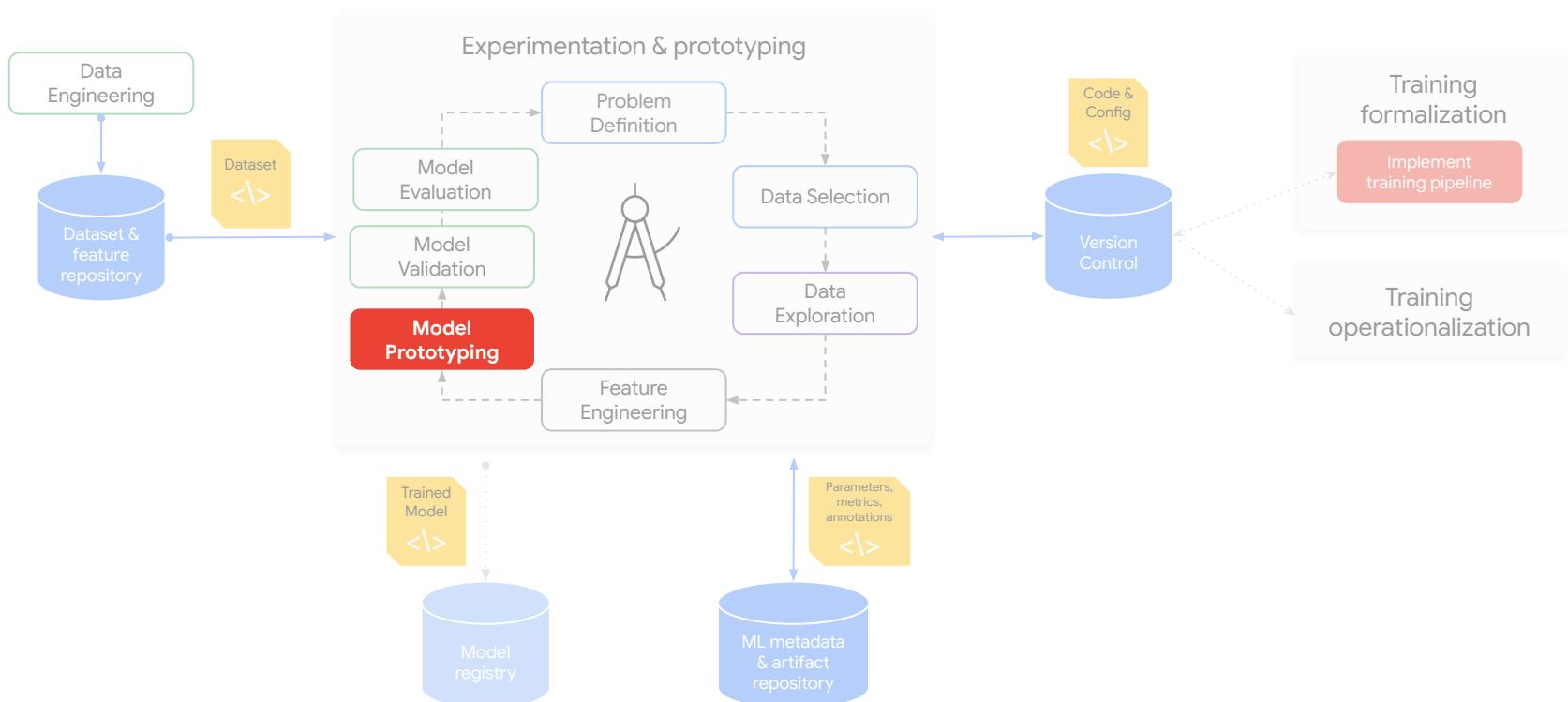


# ML Development: Model Prototyping

---

# MLOps: ML Development



# The MLOps Personas



ML  
Engineer



ML  
Researcher



Data  
Scientist



Data  
Engineer



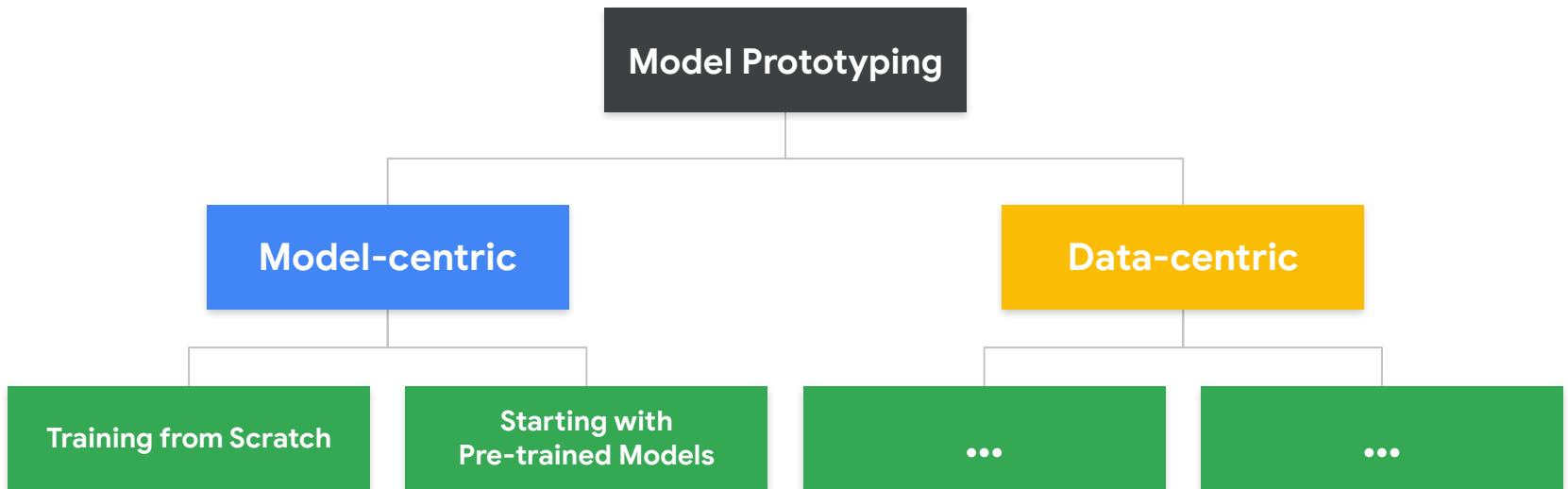
Software  
Engineer

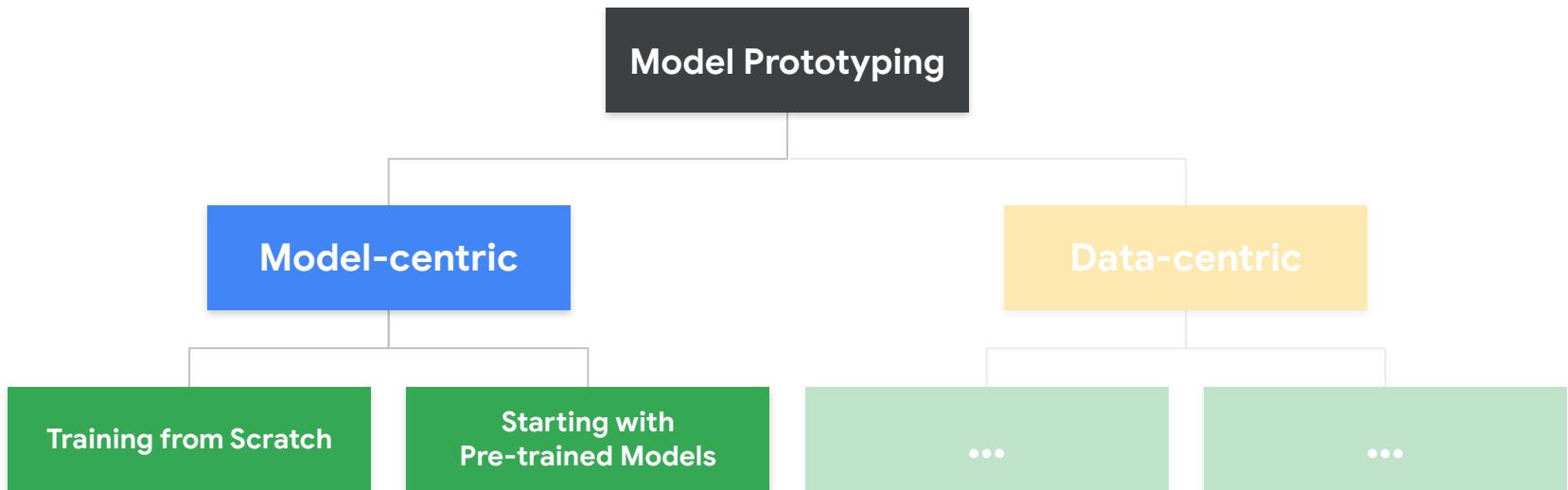


DevOps

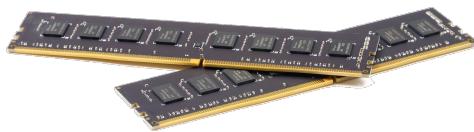


Business  
Analyst

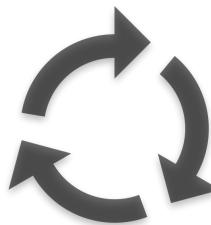




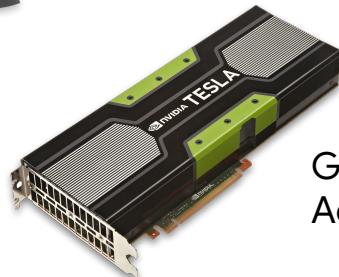
# Training Pipeline: Need Compute Resources



Memory



Compute



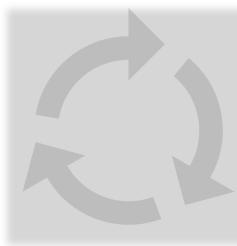
GPU and  
Accelerators

# Training Pipeline: Need Compute Resources

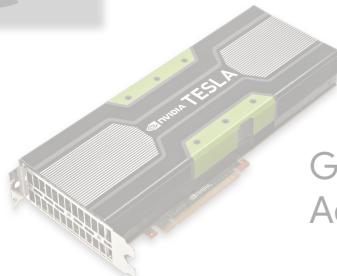
**Computationally Intensive  
Repeated Many Times (Epochs)**



Memory



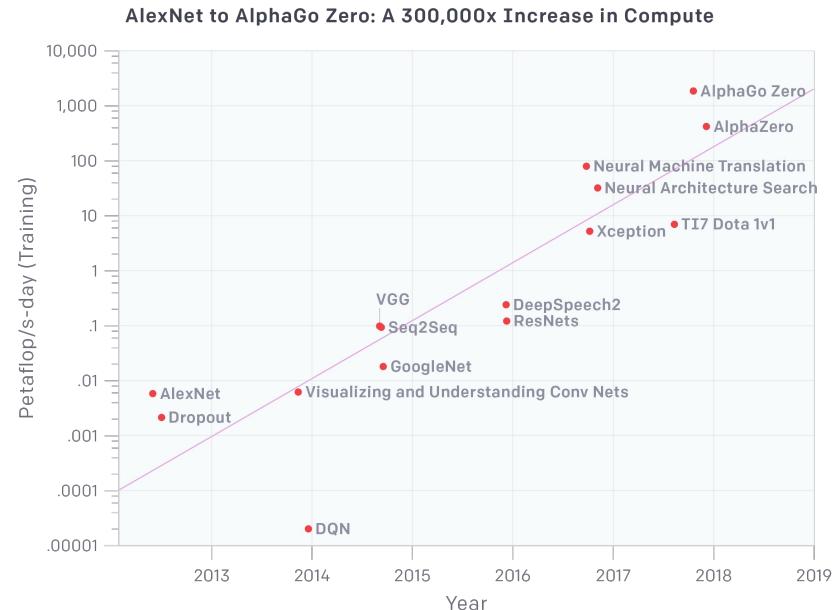
Compute



GPU and  
Accelerators

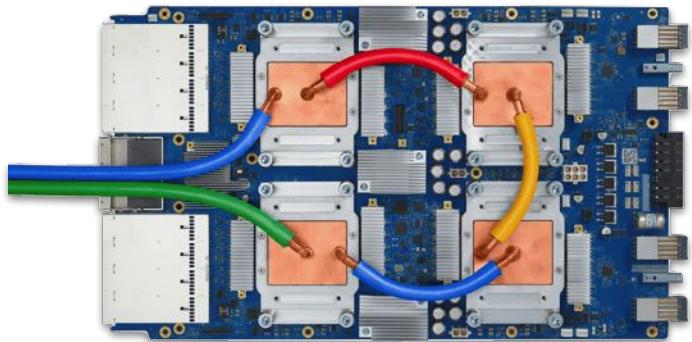
# Compute Requirements are Growing

**"... since 2012 the amount of compute used in the largest AI training runs has been increasing exponentially with a 3.5 month-doubling time** (by comparison, Moore's Law had an 18-month doubling period). Since 2012, this metric has **grown by more than 300,000x (an 18-month [Moore's Law] doubling period would yield only a 12x increase)**. Improvements in compute have been a key component of AI progress, so as long as this trend continues, it's worth preparing for the implications of systems far outside today's capabilities."



Source: <https://blog.openai.com/ai-and-compute/>

# TPUs/GPUs





Elliot Turner  
@eturner303

...

Holy crap: It costs \$245,000 to train the XLNet model  
(the one that's beating BERT on NLP tasks..512 TPU v3  
chips \* 2.5 days \* \$8 a TPU) -  
[arxiv.org/abs/1906.08237](https://arxiv.org/abs/1906.08237)

## XLNet: Generalized Autoregressive Pretraining for Language Understanding

Zhilin Yang<sup>\*1</sup>, Zihang Dai<sup>\*12</sup>, Yiming Yang<sup>1</sup>, Jaime Carbonell<sup>1</sup>,  
Ruslan Salakhutdinov<sup>1</sup>, Quoc V. Le<sup>2</sup>

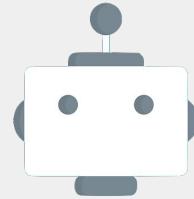
<sup>1</sup>Carnegie Mellon University, <sup>2</sup>Google Brain

{zhiliny, dzihang, yiming, jgc, rsalaku}@cs.cmu.edu, qvl@google.com

### Abstract

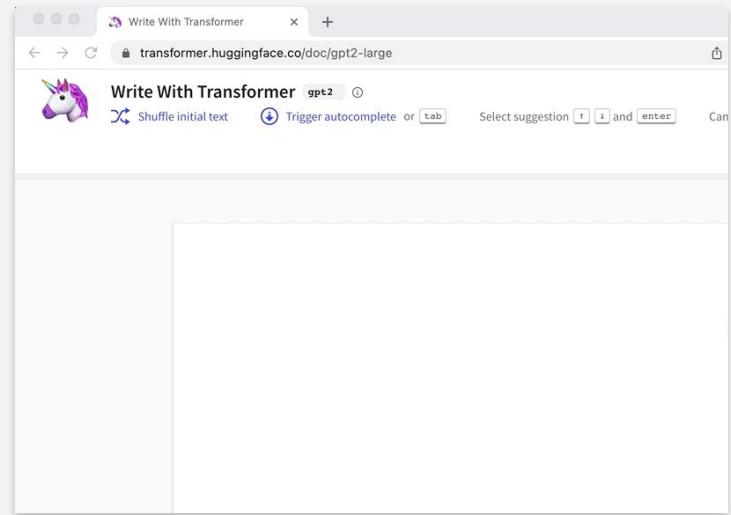
With the capability of modeling bidirectional contexts, denoising autoencoding based pretraining like BERT achieves better performance than pretraining approaches based on autoregressive language modeling. However, relying on corrupting the input with masks, BERT neglects dependency between the masked positions and suffers from a pretrain-finetune discrepancy. In light of these pros and cons, we propose XLNet, a generalized autoregressive pretraining method that (1) enables learning bidirectional contexts by maximizing the expected likelihood over all permutations of the factorization order and (2) overcomes the limitations of BERT thanks to its autoregressive formulation. Furthermore, XLNet integrates ideas from Transformer-XL, the state-of-the-art autoregressive model, into pretraining. Empirically, XLNet outperforms BERT on 20 tasks, often by a large margin, and achieves state-of-the-art results on 18 tasks including question answering, natural language inference, sentiment analysis, and document ranking.<sup>1</sup>.

Presenter HERE



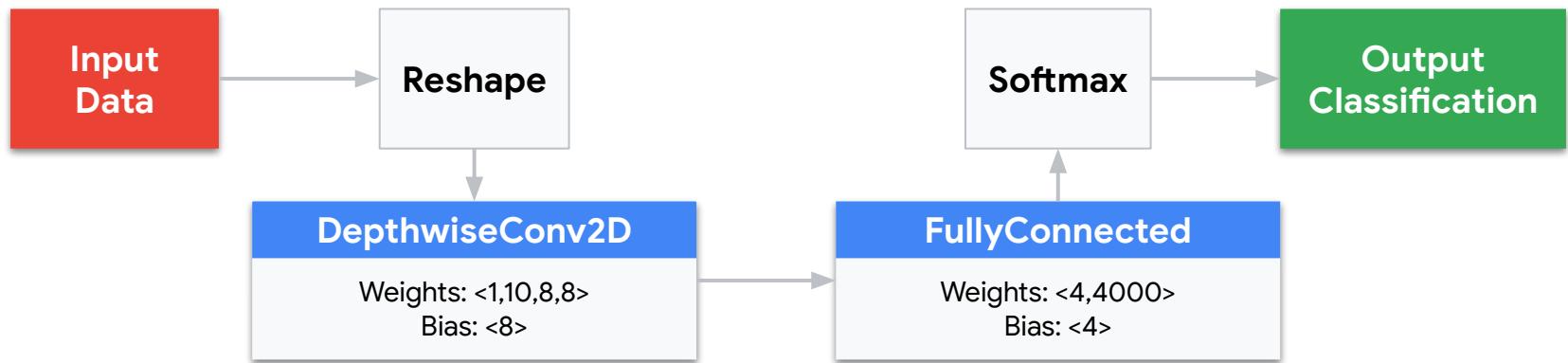
**GPT-2**

US \$256/hr



# TinyML Also Takes Long

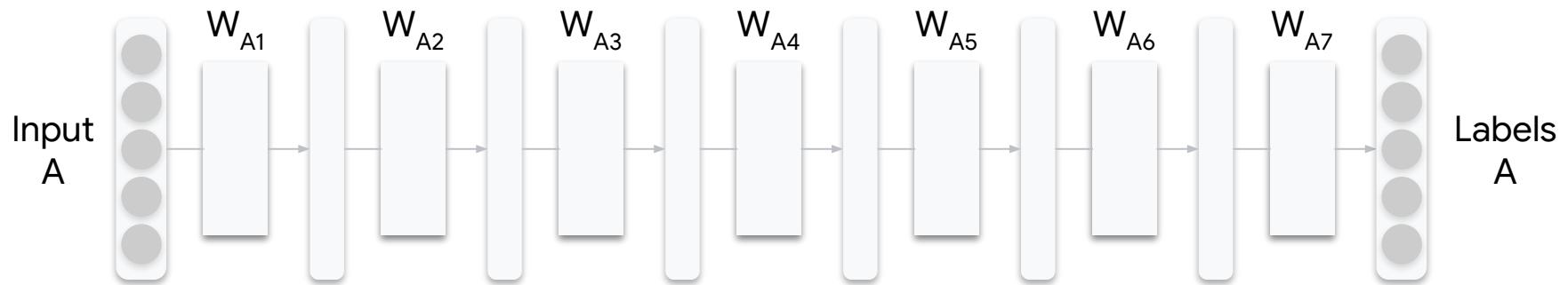
(okay, not that long...)



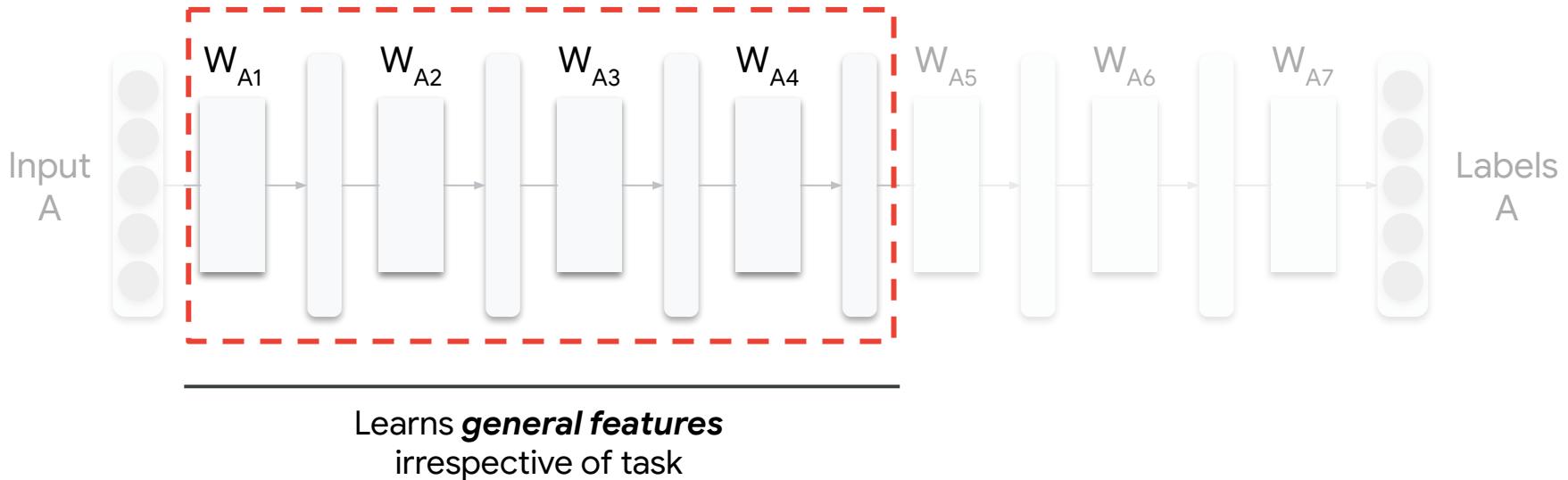
**2 hours for a 5-word model**

# Pre-trained Models

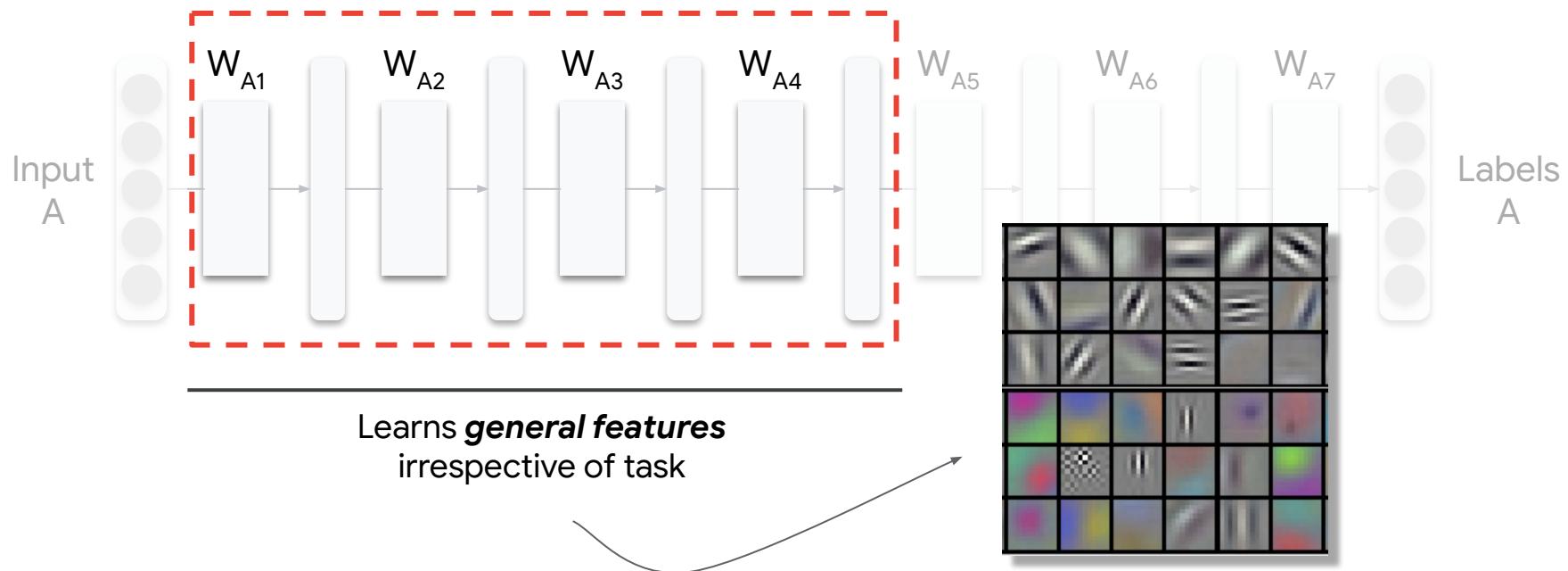
# End Result of Training



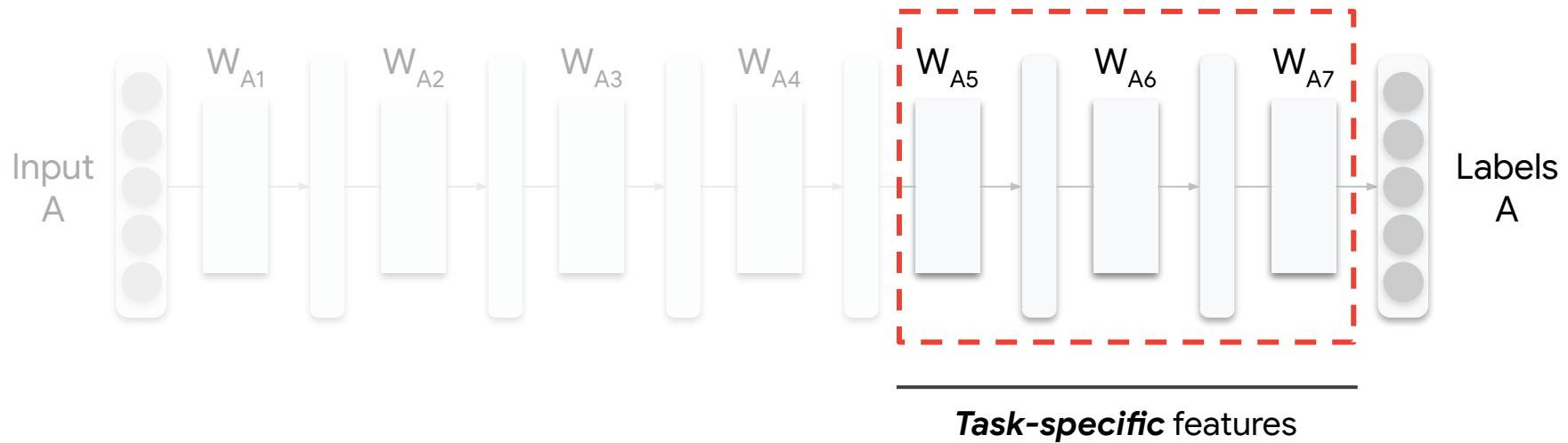
# End Result of Training



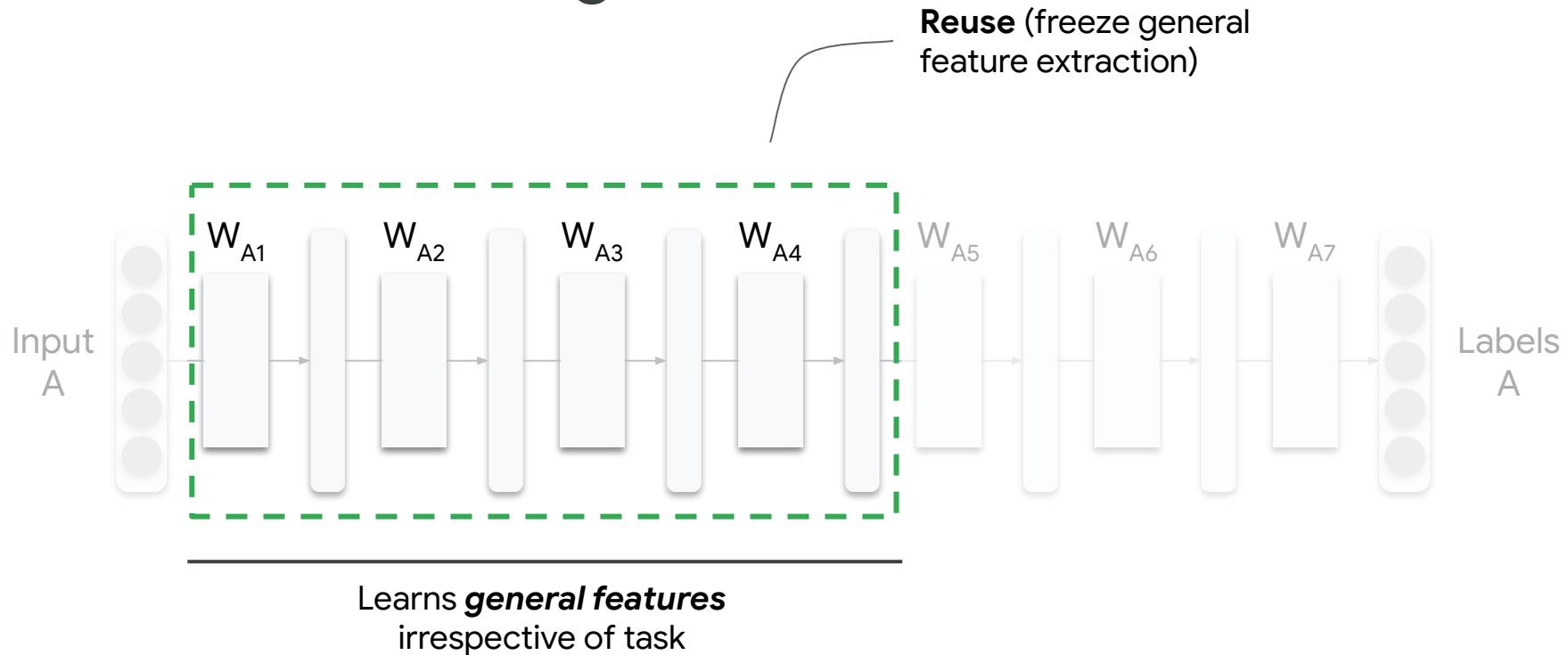
# End Result of Training



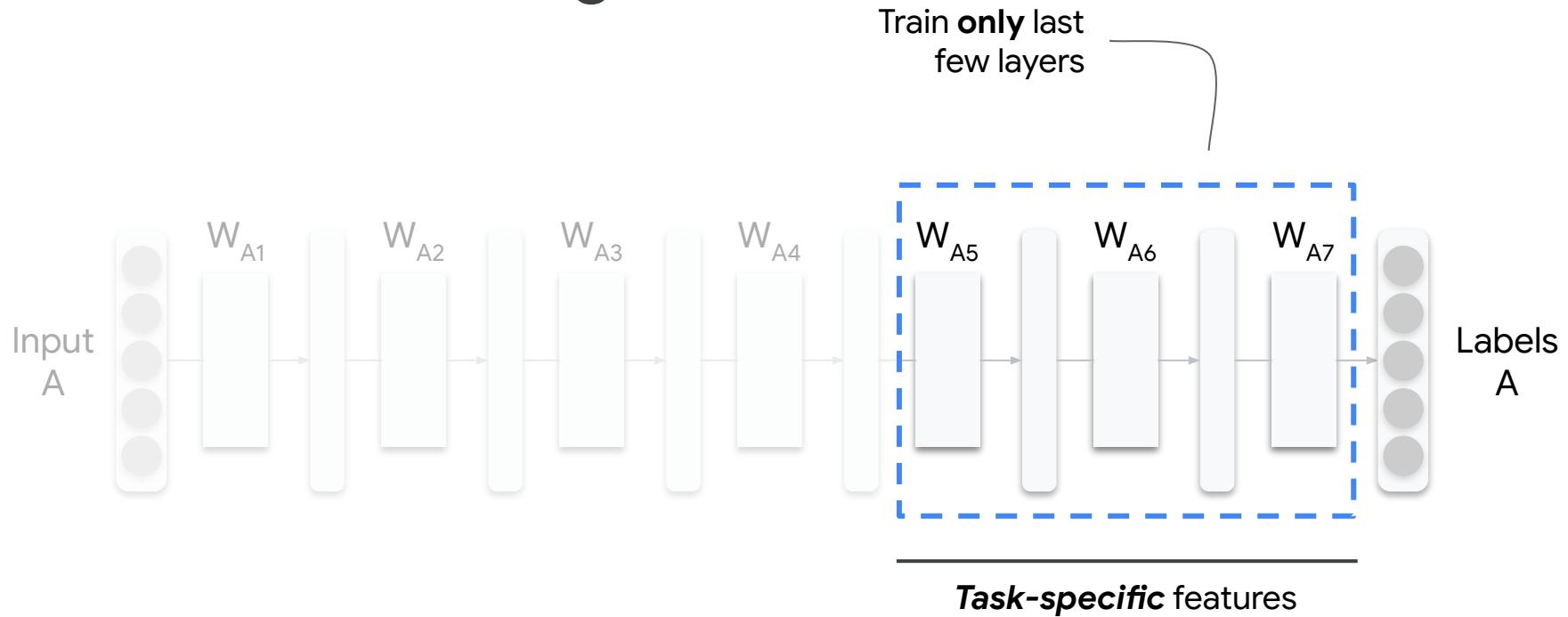
# End Result of Training



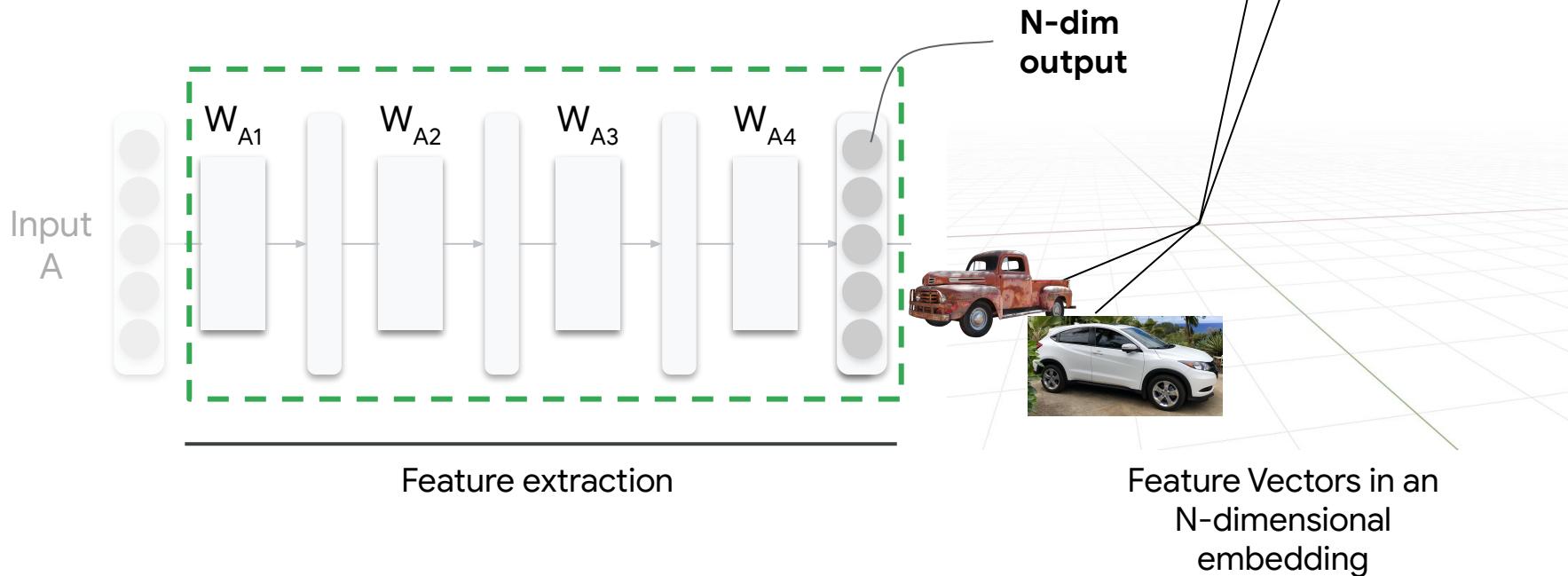
# Transfer Learning



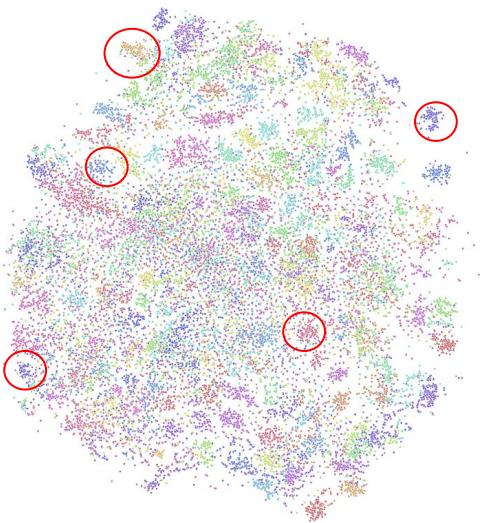
# Transfer Learning



# Embeddings

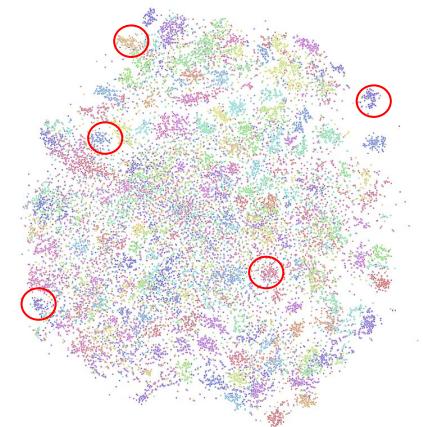
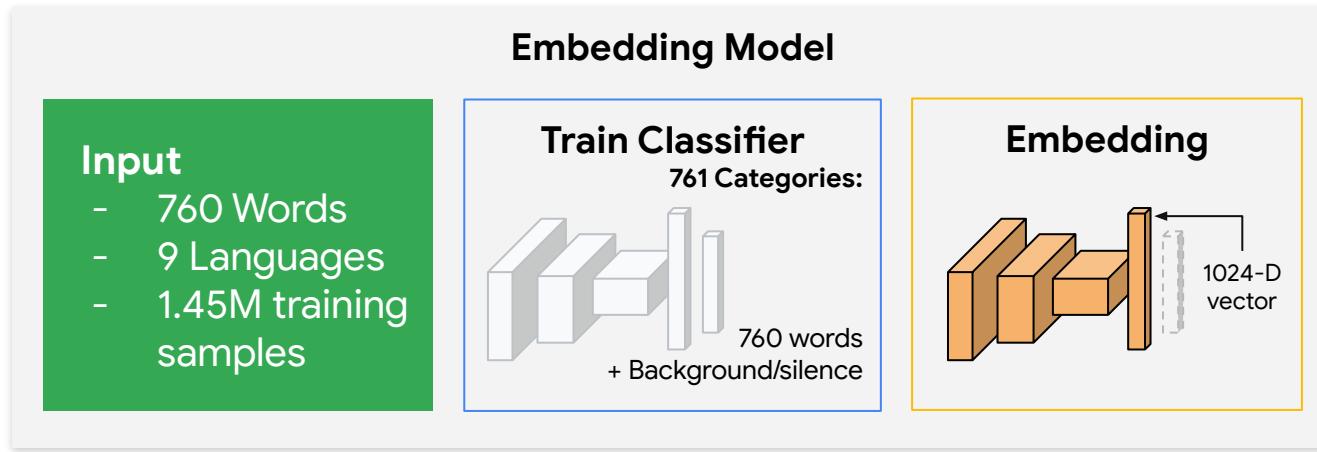


## t-SNE view of 165 novel words (not used to train the classifier)



Each cluster is a different word (each point is a  
different recording from different speakers)

# KWS Embeddings



- Minimize **intra**-class distance
- Maximize **inter**-class distance

# Few Shot Keyword Spotting in Any Language

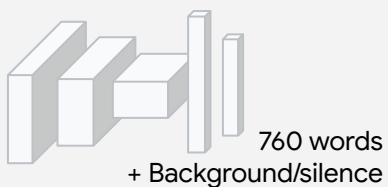
## Input

- 760 Words
- 9 Languages
- 1.45M training samples

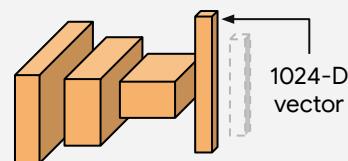
## Embedding Model

### Train Classifier

761 Categories:



## Embedding



**Input:** 5 words in any language

## Finetune KWS Model

**3 Categories**  
Target Keyword  
Unknown Word  
Background/silence

