

#### Akademia Górniczo-Hutnicza im. Stanisława Staszica w Krakowie

AGH UNIVERSITY OF SCIENCE AND TECHNOLOGY



# Hardware aware neural network compression

Krzysztof Wróbel, Marcin Pietroń, Maciej Wielgosz, **Michał Karwatowski**, Kazimierz Wiatr

Kraków 22-24.10.2018





# Natural Language Processing

- machine translation
- voice typing
- sentiment analysis
- question answering
- automatic summaryzation
- •

# Convolutional Neural Network



AGH



$$q_{\rm fxp} = \mathscr{Q}(x_{\rm flp}) = \mu + \boldsymbol{\sigma} \cdot {\rm round}(\boldsymbol{\sigma}^{-1} \cdot (x - \mu))$$

Scale factor:  $\sigma = 2^{-\text{frac}_{bits}}$  int\_bits = ceil(log<sub>2</sub>(max |x|)) Mean value  $\mu = 0$ 



www.agh.edu.pl



# Pruning

Algorithm 1 1D convolution layer with pruning

1: **for** *outS* in *output\_size* **do** 

- 2: **for** *outFM* in *output\_feature\_maps* **do**
- 3: **for** *inFM* in *input\_feature\_maps* **do**
- 4: **for** kerS in kernel\_size **do**
- 5: **if**  $abs(weight[inFM][outFM][kerS]) >= pruning_threshold$  **then** 
  - out put(outS, outFM) + = input(outS + kerS, inFM) \* weight[inFM][outFM][kerS]
  - end if
  - end for
  - end for
- 10: out put(outS, outFM) + = bias[outFM]
- 11: **end for**

12: **end for** 

6:

7:

8:

9:





### Keras model to IP core





# Movie Review dataset

- 1000 positive reviews
- 1000 negative reviews
- Validation split: 80% -20%

**POSITIVE:** the fact there's a dead body in the corner goes to enhance the feeling of paranoia and a mysterious , hurried call telling him to leave immediately is also very chilling

**NEGATIVE:** you could easily sleep through whole sections of the film ( as some fellow critics did ) and wake up in a scene exactly like the one you nodded off in , not having missed anything worthwhile





www.agh.edu.pl



www.agh.edu.pl



# **FPGA** implementation

Precision	Pruning threshold	FF	LUT
32	0	8118	60750
	0.035	4449	14906
3	0	1952	14553
	0.035	1899	14047



# Conclusions

- No more than 1% drop in accuracy after compression
- Around 4 times less FPGA resources needed



## Future work

- Non linear quantization
- Retraining
- Architectural changes





Dwa serwery z dwoma kartami firmy Nallatech zawierającymi układy FPGA Altera Stratix V

Parametry karty typ 1: Nallatech 395-AB: Stratix V AB 32GB of DDR3 Parametry karty typ 2: Nallatech 395: Stratix V D8 32GB of DDR3 Programowanie: OpenCL Altera Quartus Narzędzia: Altera OpenCL SDK Altera Quartus



