



AKADEMIA GÓRNICZO-HUTNICZA IM. STANISŁAWA STASZICA W KRAKOWIE

**DZIEDZINA: NAUKI INŻYNIERYJNO-TECHNICZNE**  
DYSCYPLINA: AUTOMATYKA, ELEKTRONIKA I ELEKTROTECHNIKA

## **ROZPRAWA DOKTORSKA**

**Sprzętowa akceleracja wymagających obliczeniowo  
operacji na potrzeby algorytmów sztucznej inteligencji  
w układach FPGA**

*Autor:* mgr inż. Michał KARWATOWSKI

*Promotor:* Prof. dr hab. inż. Kazimierz WIATR

*Promotor pomocniczy:* dr hab. inż. Maciej WIELGOSZ

*Praca wykonana:*

Akademia Górniczo-Hutnicza im. Stanisława Staszica w Krakowie,  
Wydział Informatyki, Elektroniki i Telekomunikacji, Instytut Elektroniki

Kraków, 2022



PRAGNĘ ZŁOŻYĆ PODZIĘKOWANIA DLA PANA  
PROFESORA KAZIMIERZA WIATRA ORAZ PANA  
DOKTORA MACIEJA WIELGOSZA ZA NIEOCENIONĄ  
POMOC ORAZ WSPARCIE W TRAKCIE PROWADZENIA  
BADAŃ JAK RÓWNIEŻ MOBILIZACJĘ DO PISANIA  
NINIEJSZEJ ROZPRAWY.

# Streszczenie

Niniejsza rozprawa dotyczy efektywnej implementacji algorytmów sztucznej inteligencji, w szczególności sieci neuronowych, w układach FPGA. Optymalna akceleracja sprzętowa pozwala poszerzyć zastosowania sztucznej inteligencji. Istotnym aspektem towarzyszącym badaniom prezentowanym w pracy jest optymalizacja energetyczna. Na akcelerację obliczeń składa się nie tylko dopasowana architektura sprzętowa, ale również optymalizacja algorytmów. W rozprawie badane są wpływy pruningu oraz kwantyzacji na algorytmy uczenia maszynowego, w tym sieci neuronowych. Głównymi obszarami zastosowań badanych algorytmów jest przetwarzanie języka naturalnego oraz analiza szeregów czasowych. Opracowane zostało narzędzie DL2HDL, które w istotny sposób ułatwia akcelerację w układach FPGA sieci neuronowych, zaprojektowanych w bibliotekach wysokopoziomowych, takich jak PyTorch. Konstrukcja oraz interface'y narzędzia zostały zaprojektowane w sposób przyjazny dla użytkownika, który nie posiada specjalistycznej wiedzy o projektowaniu akceleratorów FPGA. Ważnym elementem narzędzia DL2HDL jest możliwość symulacji oraz testowania działania algorytmu na każdym kroku implementacji projektu. Dzięki opracowanym optymalizacjom narzędzie generuje wydajny akcelerator zorientowany na niską latencję. Optymalizacje pod kątem latencji są wystarczająco daleko idące, aby umożliwić zastosowanie sieci neuronowych w systemach wspomagających akcelerator cząstek LHC pracujący w CERN. Cel rozprawy został osiągnięty. Badania pozwoliły na określenie optymalnych poziomów pruningu oraz kwantyzacji dla analizowanych algorytmów. Wypracowana została lista wymagań oraz powstała implementacja narzędzia DL2HDL mapującego wysokopoziomowy opis sieci neuronowych do zoptymalizowanych architektur sprzętowych w układach FPGA. Szeroko zakrojone optymalizacje pozwoliły na osiągnięcia najniższej spotkanej dotąd w literaturze latencji w trakcie inferencji rekurencyjnych sieci neuronowych typu LSTM.

# Abstract

This dissertation concerns the effective implementation of artificial intelligence algorithms, in particular neural networks, in FPGAs. Optimal hardware acceleration allows artificial intelligence algorithms to expand the field of applications. An important aspect accompanying the research presented in this dissertation is energy optimization. The acceleration of calculations consists not only of the adapted hardware architecture but also the optimization of algorithms. The dissertation examines the effects of pruning and quantization on machine learning algorithms, including neural networks. The main areas of application of the analyzed algorithms are natural language processing and time series analysis. The DL2HDL tool has been developed, which significantly facilitates the acceleration of neural networks designed in high-level libraries, such as PyTorch. The design and interfaces of the tool have been designed in the most user-friendly way, and does not require from the user any specialist knowledge about designing FPGA accelerators. An important element of the tool is the ability to simulate and test the operation of the algorithm on each step of the project implementation. Thanks to the developed optimizations, the tool generates an efficient accelerator focused on low latency. Latency optimizations go far enough to enable the use of neural networks in systems supporting the LHC particle accelerator at the CERN research center. The goal of the dissertation was successfully achieved. The research allowed to determine the optimal levels of pruning and quantization for the analyzed algorithms. A list of requirements was compiled and an implementation of a tool mapping a high-level description of neural networks to optimized hardware architectures in FPGAs was developed. Aggressive optimizations allowed to achieve state-of-the-art latency during the inference of LSTM recurrent neural networks.

# Spis treści

<b>Słownik terminów</b>	<b>8</b>
<b>1 Wstęp</b>	<b>11</b>
1.1 Motywacja . . . . .	11
1.2 Cel i tezy pracy . . . . .	14
1.3 Organizacja pracy . . . . .	16
<b>2 Wprowadzenie do tematyki</b>	<b>18</b>
2.1 Zastosowania sztucznej inteligencji . . . . .	18
2.1.1 Przetwarzanie języka naturalnego . . . . .	18
2.1.2 Analiza szeregów czasowych . . . . .	19
2.2 FPGA . . . . .	20
2.3 Obliczenia rozproszone . . . . .	23
2.3.1 Edge computing . . . . .	23
2.3.2 Map reduce . . . . .	23
2.4 Modele reprezentacji danych . . . . .	24
2.4.1 Model przestrzeni wektorowej . . . . .	24
2.4.2 TF-IDF . . . . .	25
2.4.3 Stemming . . . . .	27
2.4.4 Obliczenia rzadkie . . . . .	27
2.5 Redukcja precyzji . . . . .	28
2.5.1 Kwantyzacja . . . . .	28
2.5.2 Pruning . . . . .	30
2.6 Metody uczenia maszynowego . . . . .	30
2.6.1 K najbliższych sąsiadów . . . . .	31
2.6.2 Sieci neuronowe . . . . .	31

2.6.3	Metryki . . . . .	35
<b>3</b>	<b>Obliczenia w przestrzeni wektorowej</b>	<b>39</b>
3.1	Akcelerator sprzętowy . . . . .	39
3.2	Obliczenia rozproszone i wydajność energetyczna . . . . .	44
3.2.1	Klaster obliczeniowy z akceleratorami FPGA . . . . .	44
3.2.2	Eksperymenty . . . . .	45
3.2.3	Wyniki . . . . .	46
<b>4</b>	<b>Redukcja precyzji obliczeń - implementacje wysokopoziomowe</b>	<b>53</b>
4.1	Redukcja precyzji metryki kosinusowej w podobieństwie dokumentów . . . . .	54
4.1.1	Wprowadzenie . . . . .	54
4.1.2	Opis systemu . . . . .	54
4.2	Redukcja precyzji w algorytmie K najbliższych sąsiadów . . . . .	65
4.3	Wysokopoziomowa implementacja sieci neuronowych . . . . .	70
4.3.1	Wstęp . . . . .	70
4.3.2	Pruning . . . . .	70
4.3.3	Kwantyzacja . . . . .	73
4.3.4	Wysokopoziomowa Implementacja Sprzętowa . . . . .	73
4.3.5	Eksperymenty . . . . .	74
<b>5</b>	<b>Mapowanie sieci neuronowych do języków opisu sprzętu</b>	<b>78</b>
5.1	Wprowadzenie . . . . .	78
5.2	Implementacja . . . . .	83
5.2.1	Wytyczne . . . . .	83
5.2.2	Ogólny opis całości procesu generacji modelu . . . . .	86
5.2.3	Elementy łączące . . . . .	91
5.2.4	Warstwy . . . . .	94
5.2.5	Inne optymalizacje . . . . .	114
5.2.6	Symulacja i uruchomienie kodu . . . . .	116
5.3	Eksperymenty . . . . .	119
5.3.1	Pasażerowie międzynarodowych linii lotniczych . . . . .	119
5.3.2	Zanieczyszczenie powietrza w Pekinie . . . . .	123
5.3.3	Magnesy nadprzewodzące . . . . .	125

5.3.4	Porównanie z innymi narzędziami . . . . .	129
5.4	Podsumowanie . . . . .	132
<b>6</b>	<b>Podsumowanie</b>	<b>133</b>
	<b>Bibliografia</b>	<b>138</b>
	<b>Załączniki</b>	<b>146</b>