

Artificial Intelligence for Mobile Communications

Dr. Dirk Wübben

25. Fachtagung Mobilkommunikation 3.-4.11.2021

 $C^{P} = ONC$ ip cores & system solutions





Deutsches

für Künstliche

Intelligenz GmbH

Forschungszentrum



Bremen







GEFÖRDERT VOM

*

Bundesministerium für Bildung und Forschung

Funkkommunikation mit Künstlicher Intelligenz (FunKI)

Radio Communication with Artificial Intelligence



- Volume
 - 6.28 M€
- Duration
 - 05/2020 05/2023
- Coordinator
 - University of Bremen Prof. Dr. Armin Dekorsy Dr. Dirk Wübben
- Internet
 - www.funki.tech

GEFÖRDERT VOM





Problem Statement

- Modern wireless communication systems are required that can support massive increase of devices while being both powerful and resource-efficient
- Higher protocol layers have successfully applied methods of artificial intelligence (AI)
- Model-driven approach for PHY and MAC are common
 - Good and practicable models used under idealized assumptions
 - Models enable the development of efficient procedures and algorithms
- Challenge: Implementations with traditional methods become increasingly difficult
 - Conflicting requirements in terms of reliability, transmission rate, number of terminals per area, latency, available radio resources, energy efficiency, complexity, and hardware costs
 - Mathematical modelling of the problem is no longer possible (model deficit) or leads to very complex models (algorithm deficit)
- Data-driven machine learning (ML) methods represent a promising approach for gaining an understanding of the model or for system and technology design



Source: http://www.emfexplained.info

Gb/s pJ/bit 10⁻⁹ 10000/km² < 1 ms mm² €





Scopes of Learning

- Offline Learning
 - Replacement of *classical* signal processing tasks by pre-trained ML-blocks (e.g. NNs)
 - Training is done offline and resulting ML-block is finally used in the actual system
 - Properties: May reduce latency and/or improves the resulting BER performance

Online Learning

- Implementing AI at run-time to adapt/refine a communication system to a specific channel or to hardware impairment
- On the fly adaption to effects that are unfeasible to model in classical channel models
- Example: Adaptive end-to-end training of the autoencoder with continuously re-training

Design-time Learning

- Transfer insights from the learning process back to classical signal processing algorithms by applying the outcome in classical receiver structures
- Leverage existing transceivers by means of machine learning and by new insights from learning algorithms



AL CONTRACTOR

NOKIA

ReefShark

Source: https://macropolo.o







Main topics

Conceptual design, optimization and evaluation of essential components

- Transmitter, channel, receiver
- Channel estimation, prediction, feedback

FPGA hardware implementation / ASIC design

• Implementation of NN-based transceiver components by hardware accelerators



Training data generation

- Measurement prototype for multi-antenna/multi-carrier channel measurements
- Measured/real channel data for massive MIMO systems

Testing and demonstration of the technologies in selected fields of application

• SDR-based demonstrator implementing selected AI components



Workpackages



Scenarios, KPIs
Exchange of trained

- transceivers and of channel measurements
- Iterative workflow for the transceiver implementation ("bidirectional research approach")



5 Exchange of knowledge and results



WP2 Transmitter and receiver structures



- Motivation
 - SotA: Handcrafted processing blocks derived from channel models, e.g., coding, pulse shaping, channel estimation
 - Those models can only generalize or approximate the real channel
 - \rightarrow mismatch between processing block and the real channel possible
 - Dividing a system into single and separate blocks may not be optimal
 - How to combine several communication blocks?
 - No single model can cover all realizations
- Goal
 - Optimization of single components of as well as the whole communication system by application of deep learning

Approaches

- Measurements and a following data-driven optimization of single blocks
- Redesign of a communication system to allow endto-end training (e.g. autoencoder)





WP3 Channel adaptation

- A very well tuned channel adaptation helps to realize the gains promised by enhanced algorithms (MU-MIMO, Beamforming, CoMP, CoSch, etc.)
- Parameter estimation and CSI prediction is the basis for channel adaptation
- PRACH measures about frequency and timing offset is key for further parameter estimation and prediction
- Goal
 - Develop AI based algorithms that enhance system behavior by improvement of PRACH measurements, channel estimation & prediction, and channel adaptation

Approaches

- Channel measurements in relevant scenarios and according channel modelling
- Al based profiling of multipath components for channel estimation & prediction
- AI based PRACH detection and parameter estimation with expert knowledge
- AI based channel adaptation that adapts to different situations and environments



Link quality for sophisticated algorithm at reception time?

multipath components evolving over time

MOTIUS





Interference at reception time?

WP4 Efficient transceiver implementation

Motivation

- Baseband signal processing components, e.g. demodulation, equalization, channel decoding, are extremely critical building blocks w.r.t. throughput, latency and energy efficiency in the overall system and must thus be implemented as dedicated hardware accelerators
- Large gap between "high-level" complexity analyses and actual implementation complexity

Goal

 Implementation and evaluation of AI processing elements as dedicated hardware accelerators in ASIC and FPGA technologies and fair comparison with state-of-the-art solutions

Approaches

Funkl

- Implementation/simulation frameworks for fast design space exploration
- Hardware/algorithm co-design
- Investigate trade-off between implementation cost and performance





ASIC implementation 22nm FD-SOI: 4 and 8 iteration LDPC channel decoder



FPGA implementation: Autoencoder with trainable demapper



WP5 Testbed implementations



Motivation

Various ideas and scenarios are introduced in other WP. But are they doable in existing hardware?

Goal

- SDR based Testbed for implementation and evaluation of ML-based baseband processing
- Show the feasibility to implement the various algorithms in specific scenarios
- Comparison with existing solutions to demonstrate possible advantages of KI in communications

Approach

- Mixed Architecture: selected KI component on FPGA (for example Xilinx Alveo), remaining test chain in GNU-Radio
- 5G Open-Source Software







ML-based Decoder Implementations



Forward Error Correction





- Forward Error Correction (FEC)
 - Adding redundancy to improve reliability of message exchange
- Problem Statement
 - Decoder is the most complex part of the signal processing chain
 - High energy consumption / chip area by decoder implementation, iterative decoder introduces additional latency
- Goal
 - Decoder implementations meeting trade-off wrt. energy-efficiency, low-complexity, e2e performance, and latency
- Approaches: Designing decoder schemes using machine-learning driven approach
 - MLS-BP: Adapt common decoder implementation for short block length
 - Discrete Decoder: Implement decoder with very low bit resolution
 - NN-FoC: Execute decoder only if success is likely



Decoding of LDPC-Codes via Belief Propagation (BP)





- Message Passing (MP) / Belief Propagation (BP) [KFL01]
 - Decoding by exchanging messages between variable nodes (VN) and check nodes (CN)

Initialize VN messages





Unrolled Factor Graph

- Ig factor graph to visualize the message propagation over iterations \Rightarrow favorable from implementations \Rightarrow favorable from the interactions \Rightarrow favorable from the interactio
- Unrolling factor graph to visualize the message propagation over iterations \rightarrow favorable HW implement.

- Problem: BP decoder is suboptimal for short block length
 - Performance degrades due to loops in the factor graph ightarrow reliability of messages are overestimated
- Neuronal BP (N-BP) interprets graph as Neural Network with adaptable weights [NBB18]
 - Training of N-BP starts as standard BP decoder
 - Drawbacks: training for fixed number of iterations, huge training complexity (increase in k), number of weights scale with iterations



ML-Scaled Belief Propagation for Short Codes



ML-Scaled Belief Propagation (MLS-BP)





- Alternative Approach
 - Use single scalar factor $\beta \leq 1$ for received LLRs and common BP \rightarrow equivalent channel reliability factor $L_{\text{MLS}} = \beta L_{\text{ch}}$
- Bruteforce (BF) approach to determine L_{MLS}
 - Per SNR loop over scaling factors, simulate BER and select $L_{\rm BF}^*$ leading to minimum BER
- Machine Learning Scaled Belief Propagation (MLS-BP) [HWD21]
 - Adaptation of L_{MLS} via simple ML approach
 - Training: Generate training dataset (y, x) and minimize loss function MSE loss
 MI loss

$$\mathcal{L}_{\text{MSE}} = \|\boldsymbol{x} - \hat{\boldsymbol{x}}\|_2^2$$

$$\mathcal{L}_{\mathrm{MI}} = -\left(1 - \frac{1}{n} \sum_{i=1}^{n} \log_2\left(1 + e^{-x_i \cdot L(x_i)}\right)\right)$$

Inference: No extra overhead during execution





[HWD21] Hummert, Wübben, Dekorsy: Machine Learning Scaled Belief Propagation for Short Codes, VTC-Fall 2021

Performance for (63, 45) BCH code

Scenario

(63, 45) BCH code, 5 BP iterations, BPSK, AWGN

Observations

- N-BP shows substantial gains to standard BP (loss due to loops in the graph)
- MLS-BP is nearly indistinguishable to N-BP
- Range-learned constant L^{range}_{MLS}: Use training set containing a range of SNR to learn scaling factor
 → nearly same performance as SNR-specific L_{MLS}
- As expected, L^{*}_{BF} shows same performance

Conclusion

Funkl

 Range-learned constant is sufficient to improve the performance of BP decoder for short codes





Observations

LDPC codes for different length, 5 BP iterations and

- With increasing codeword length n, suboptimality of BP reduces
- Possible gains of MLS-BP reduce as well

Conclusion

Funkl

Scenario

- MLS-BP introduces trainable scaling factor for LLRs
- Instead of SNR-specific scaling factor one rangedlearned factor is sufficient
- No adaption of BP implementation necessary
- For short codes performance improvements are achieved *for free*

Performance for longer LDPC codes

range-learned constant $L_{\rm MIS}^{\rm range}$





Discrete Decoder Implementation



Information Bottleneck Method for Discrete Decoder Design



- Design of efficient decoding algorithms with low bit-resolution
- Approach
 - Learn information processing chain by focusing only on relevant information while maximizing transmission rate
 - Mutual information to formulate trade-off between complexity and relevant information *I*(x; z)
- Information Bottleneck Method (IBM)
 - Relevant information processing: interest is on information of source signal x or u
 - Quantization: Trade-off between compression rate I(y; z) and relevant information I(x; z)
 - IB-based decoder: Low internal variable resolution (e.g. 3-4 bit) and simple discrete operations



Universität

Bremen

LUT-MP: Decoder Learning





• LUT-MP decoder algorithm is learned via IB ($\beta \rightarrow \infty$) in **Discrete Density Evolution (DDE)** [KYK08, KY14]





Kurkoski, Yamaguchi, Kobayashi: Noise thresholds for discrete LDPC decoding mappings, IEEE Globecom, 2008 Kurkoski, Yagi: Quantization of Binary-Input Discrete Memoryless Channels, in IEEE Transactions on Information Theory, Aug. 2014

LUT-MP: Decoder Execution





- Lookup-Table based Message Passing (LUT-MP) [RK16]
 - Node operations in variable nodes (VN) and check nodes (CN) are implemented by fixed LUTs



Decision $\hat{c}_{j}^{(i)} = \text{LUT}_{\text{DC}}^{(i)} \left(\bar{t}_{j \to m}^{(i)}, t_{m \to j}^{(i)} \right)$ Initialize VN messages

 $\overline{t}_{j \to m}^{(0)} = \overline{z}_j \ j = 1, \dots, 8 \ m \in \overline{\Omega}_j$





Min-LUT [MM20]: minimum of integers

$$LUT_{MS}(\bar{t}) = \begin{cases} \min_{k} |\bar{t}_{k}| & \text{if } \prod_{k} \operatorname{sign}(\bar{t}_{k}) = 1 \\ \neg \min_{k} |\bar{t}_{k}| & \text{if } \prod_{k} \operatorname{sign}(\bar{t}_{k}) = -1 \end{cases}$$



[RK16] Romero, Kurkoski: LDPC Decoding Mappings That Maximize Mutual Information, IEEE JSAC, 2016
[MM20] Meidlinger, Matz, Burg: Design and Decoding of Irregular LDPC Codes Based on Discrete Message Passing, IEEE TCOM, 2020

Performance Evaluation



Simulation setup

- (3,6)-regular LDPC Code, n = 816, $i_{max} = 40$
- Benchmark: **FP-SPA** with double precision
- Observations
 - FP-SPA with 3-bit IB ADC → 0.1 dB loss
 - FP-MS: Min-Sum with 3-bit IB ADC → 0.45 dB loss
 - 3-bit LUT-MP → 0.15 dB loss
 - **3-bit Min-LUT** → almost no difference
- Conclusion
 - Discrete decoder with 3-bit resolution performs close to SPA with double precision
 - 4-bit LUT-MP would reduce gap to FP-SPA





LDPC Decoder – First synthesis results

- Example code: n = 24, k = 12, WiMAX-like code, $d_v \in \{2, 3, 6\}$ and $d_c = 6$
- Synthesis in 28nm FD-SOI technology
- Comparison 6-bit Min-Sum vs. 3-bit LUT-MP
 - Reduction of cell area by 14%
 - 50% reduced node interconnections due to reduction from 6 to 3 bit
- Comparison 6-bit Min-Sum vs. 3-bit Min-LUT
 - Reduction of cell area by 38%
- Comparison 6-bit Min-Sum vs. 4-bit LUT-MP
 - Larger total cell area, but 33% reduced node interconnections (6 to 4 bit)

Architecture	MS 6bit	LUT-MP 4bit	Min-LUT 4bit	LUT-MP 3bit	Min-LUT 3bit
Codeblock Size	24	24	24	24	24
Pipeline Stages	10	10	10	10	10
Target (ns)	2.00	2.00	2.00	2.00	2.00
Critial Path (ns)	1.75	1.79	1.49	1.49	1.51
Frequency (MHz)	571	559	500	671	662
Throughput (Gbps)	13.7	13.4	12.0	16.1	15.9
Total Cell Area (umm ²)	55711.9	143553.9	95262.9	47908 6	34588.6
Area wrt. MS 6bit	100%	257%	171%	86%	62%







Relative Entropy based Message Combining (REMC)





- Problem: LUT-MP has been designed for specific e2e distribution
- Challenge: How to apply given LUT-MP implementation for other scenarios, e.g. higher modulation order, diversity, discrete signals with different cardinality
- Relative Entropy based Message Combining (REMC) [MWD21]
 - Combine messages with similar meaning $p(c_v | z_1, ..., z_J)$ wrt. given decoder design distribution $p^*(c|t)$

$$\begin{aligned} z_{\nu} &= \mathbf{Q}_{C,\nu}(z_1, z_2, \dots, z_J) \\ &= \operatorname*{argmin}_{t \in \mathcal{T}} D_{\mathrm{KL}}(p(c_{\nu} | z_1, z_2, \dots, z_J) || p^*(c|t)) \end{aligned}$$

Forwarding 3-bit messages + LUT-MP decoder close to optimum





NN-based Forecasting of Decodability for Early ARQ





- **Problem**: Iterative decoder introduces complexity and additional latency
 - ARQ asks for re-transmissions in case of decoder failure
- Idea: Forecast decoder success based on received LLR Ly without executing the decoder
 - Save computing resources by not running complex decoder if not promising
 - Early ARQ: Feedback of (forecasted) NACK immediately initiates re-transmission → reducing latency (e.g., URLLC)
- NN-based Forecast (NN-FoC) of Decodability [HWD21]

Funkl

- Use NN to determine probability $p(s|L_y)$ of decoding packet correctly (s = 1) for given receive LLR L_y
- Hard decision with threshold α to decide $\hat{s} = Q_{\alpha} \left(p(s|L_y) \right)$
- Structure: Feedforward NN with input layer of length *n*, hidden layers with ReLu activation, scalar sigmoid output
- Training dataset with features L_y and labels $s \in \{0,1\} \rightarrow (L_y, s)$

Performance Evaluation for different NN-FoC



System

- False Forecasts $P_{\text{FF}} = \frac{\# \text{ of false forecasts}}{\text{ all forecasts}} \text{ vs SNR}$
- (7, 4) Hamming code and BP decoder
- NN-FoC with high, medium, and low complexity
- SNR-based classifier that forecast all packets to be decodable → P_{FF} = 1 − FER
- Benchmark: High complex NN-FoC trained per SNR

Observations

- NN-FoC trained per SNR outperforms FER-classifier
- All NN-FoC with range training perform almost as good as the NN-FoC with training per SNR
- Conclusion
 - Low complex NN-FoC with training over SNR range is sufficient

NN-FoC	High compl.	Mid compl.	Low compl.	
Hidden layer	4	4	2	
Width of layer	200, 100, 500, 200	50, 50, 50, 20	50,20	
# of weights	2.617.901	9.091	1.441	
10^{0} 10^{-1} 10^{-2}		NN-FoC, Hig NN-FoC, Me NN-FoC, Lo NN-FoC 1, P 1-FER	gh compl. edium compl. w compl. Per SNR	
10 ⁻³				
10^{-4}) 1 2 3	4 5	6 7	



Efficiency Analysis



• Finite state diagram to analyze efficiency η of NN-FoC with E-ARQ and different decoder delays $\kappa = T_{Dec}/T_B$



$$\eta = R_c \frac{T_B}{T_{avg}} = R_c \frac{(1 - (1 - P_s) - P_{\rm FP}P_s)^2}{P_{\rm TP}P_s ((\kappa + 1)P_s + (1 - P_s))}$$

System

- For (32, 16) LDPC code comparison of common ARQ and E-ARQ with genie classifier and NN-FoC
- Conclusion
 - Efficiency improvements by NN-FoC and E-ARQ in comparison to standard ARQ, gains increase for higher κ





Conclusion

• Challenge: FEC decoder is most complex component of receiver chain

Machine Learning Scaled Belief Propagation (MLS-BP)

- Limit impact of overestimating LLRs for short block length by one trainable scaling factor
- For short codes performance improvements are achieved for free

Discrete Decoder Design

- Information Bottleneck Methods provides approach to learn information processing with minimum complexity
- LUT-MP decoder with 3-bit achieves performance close to optimum
- **REMQ** combines messages with similar meaning, e.g. to apply fixed decoder for general systems
- NN-based Forecast (NN-FoC)
 - Forecasting of decodability to save decoding complexity and latency

