



***Integrated Circuits for Communications
after the Happy Scaling Era***
(based on examples from channel coding)

Andreas Burg

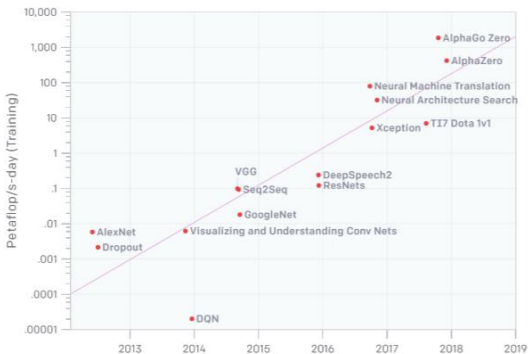
Telecommunications Circuits Laboratory, EPFL



The End of the "Happy-Scaling" Era: Lets get to work!



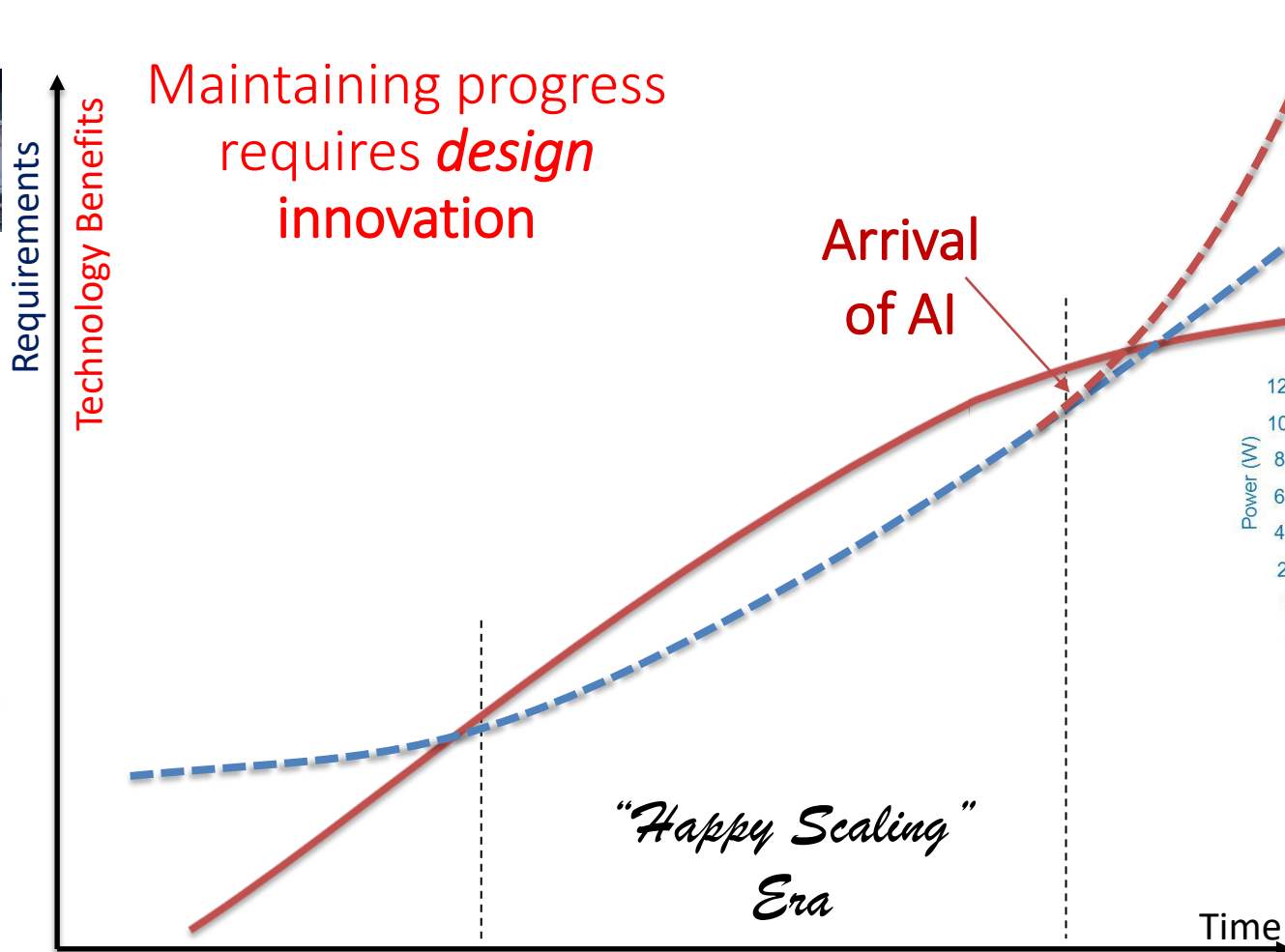
New Applications: AI



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



More processing in a reduced power envelope



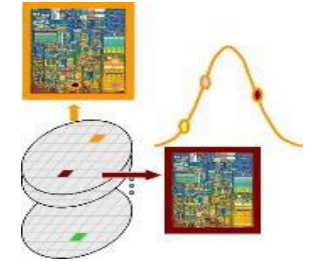
90 nm



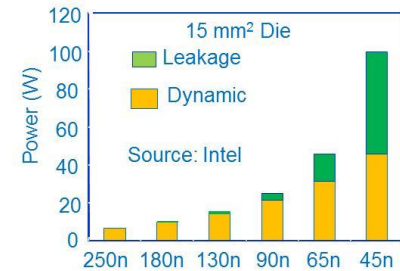
28 nm

TODAY

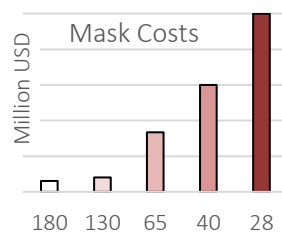
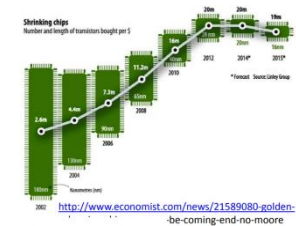
Process



Severe variability and reliability issues



Vanishing energy and performance benefits

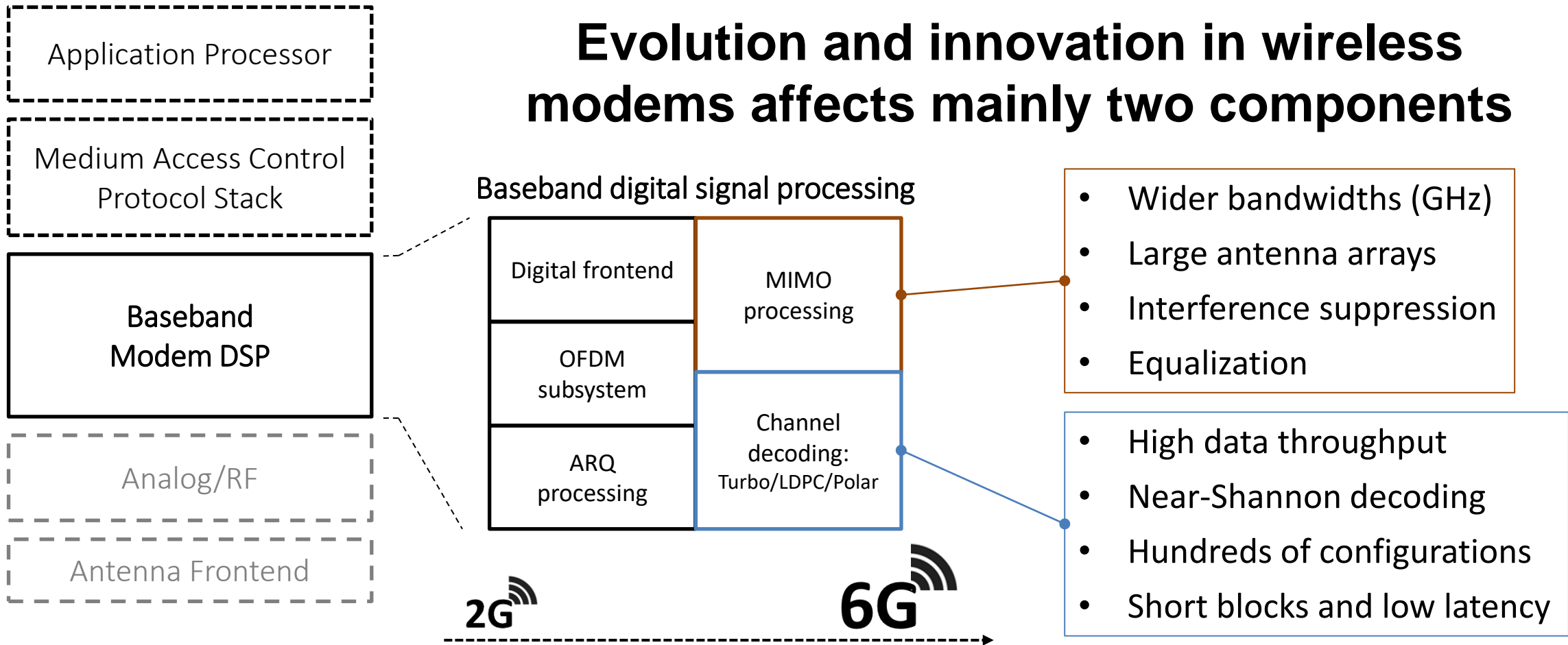


Skyrocketing costs (mask-, silicon, and design-costs)

Wireless Modem Architecture & Trends

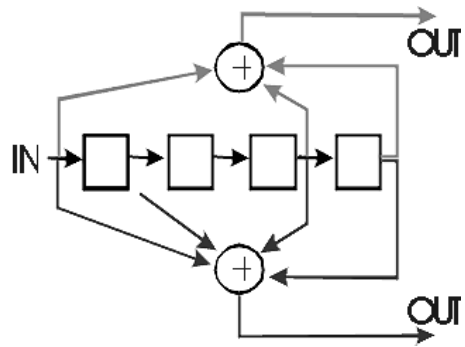
Wireless modem is a key component in mobile SoCs

- Affected by constant updates and innovations as algorithms and standards evolve

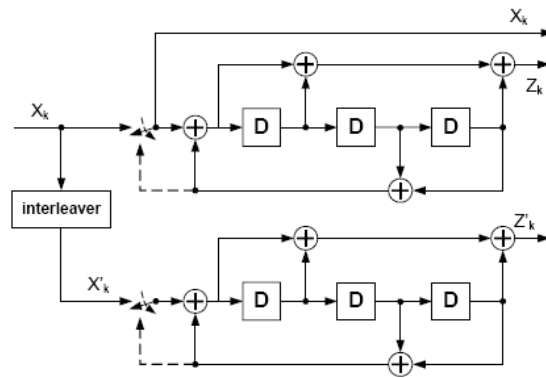


Channel Codes in 3GPP Mobile Wireless Systems

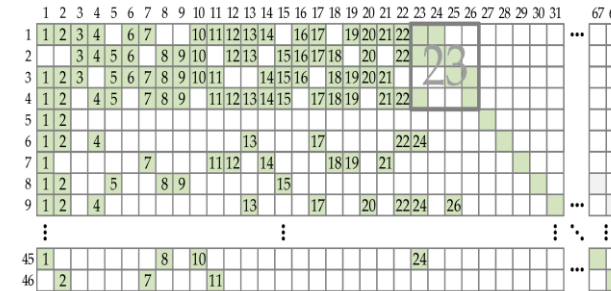
Convolutional Codes



Turbo Codes



LDPC / Polar Codes



2G 

3G 

4G 

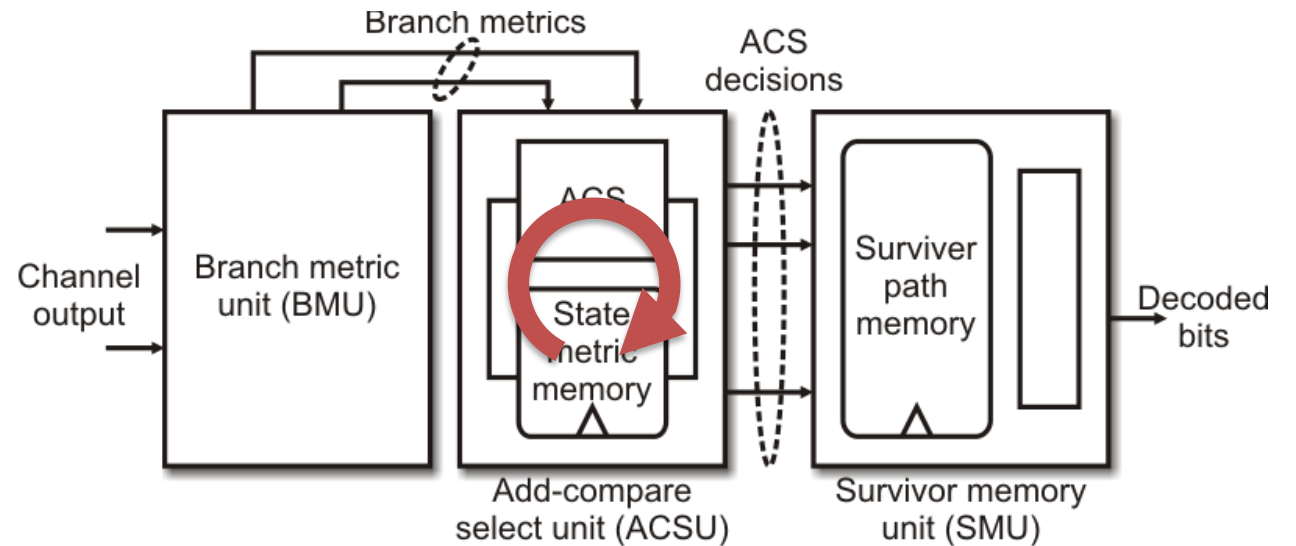
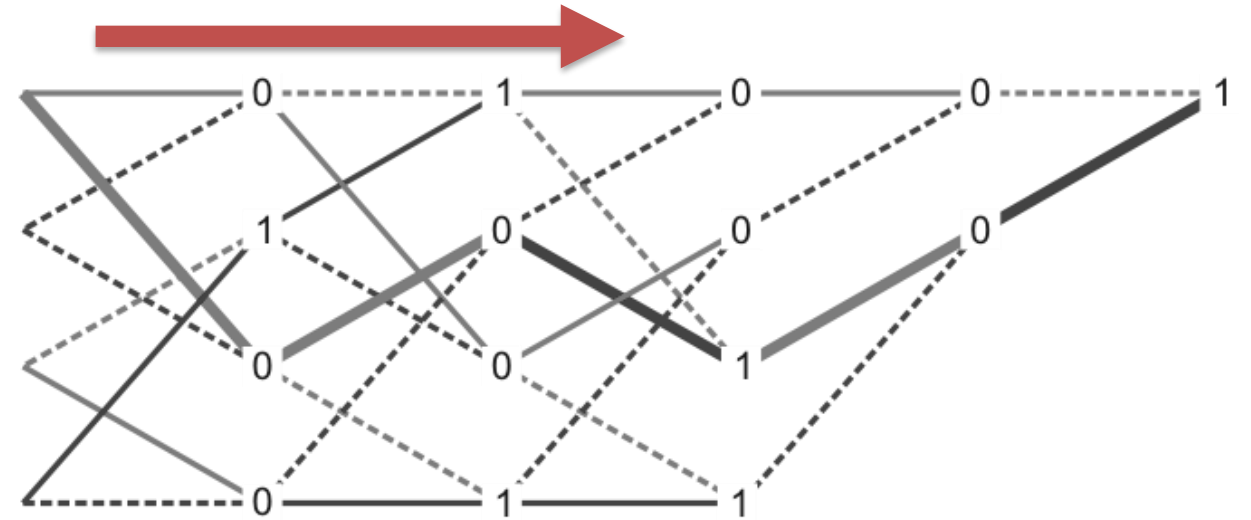
5G 

6G 

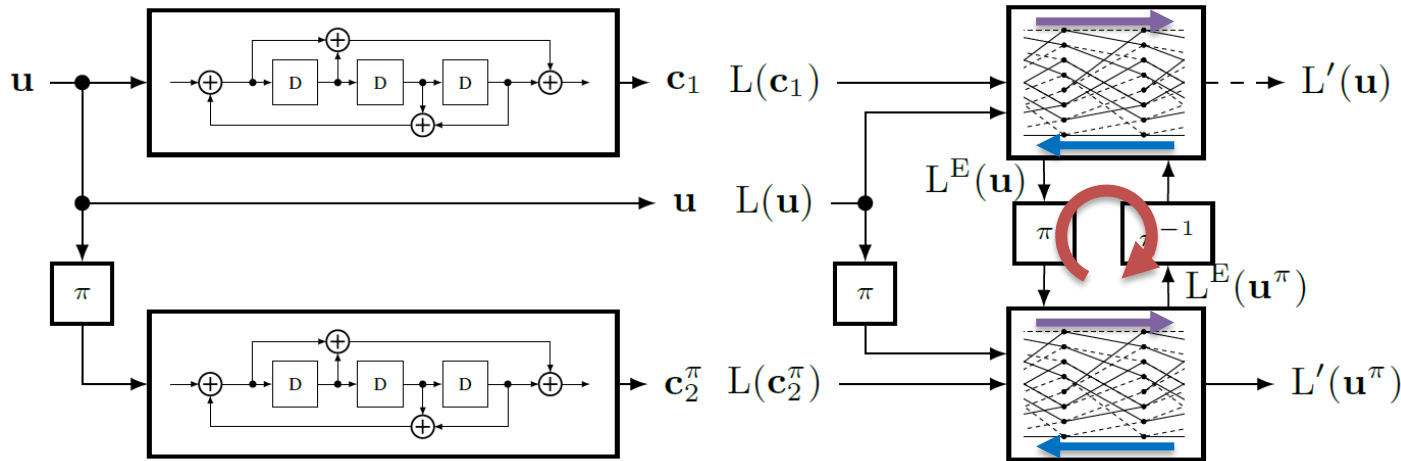
Sequential Decoding of Convolutional Codes

VITERBI algorithm (1967)

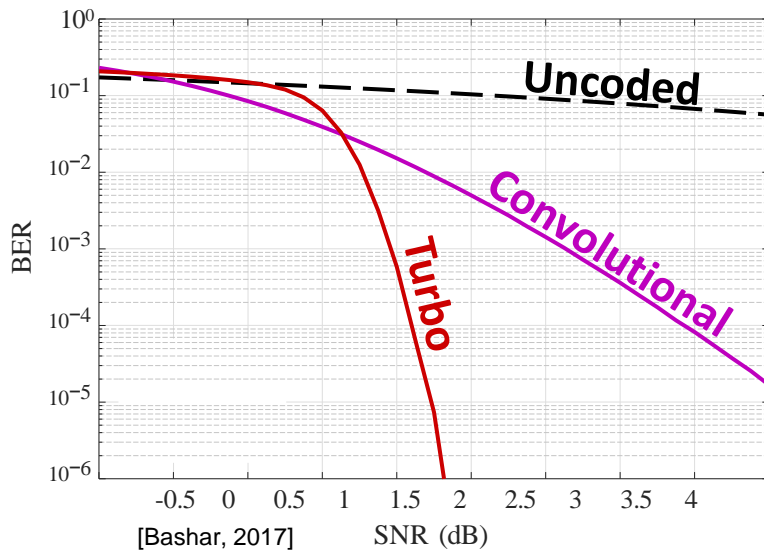
- Step-by-step trellis traversal
- Inherently **sequential process**
 - Bits per clock cycle defined only by the code
 - **Parallel processing only with exponential complexity**



3G: From Convolutional to Turbo Codes (1990s)



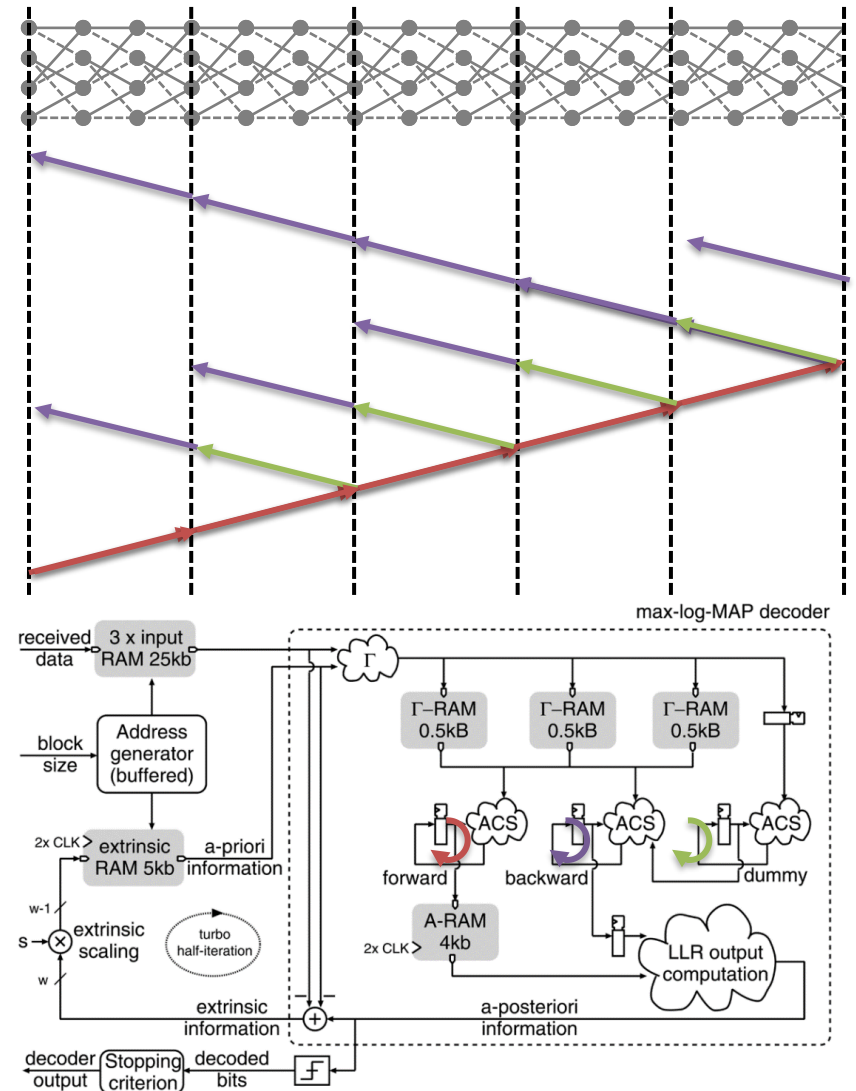
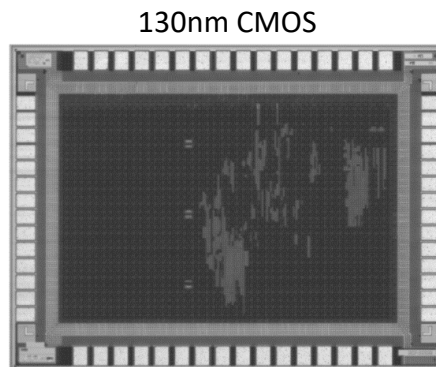
- Encoding with two independent convolutional codes based on “randomly interleaved” data
- Iterative decoding
 - Two component decoders operating on interleaved convolutional codes
 - Exchange of reliability information
 - Component decoders based on the BCJR algorithm
- Performance close to Shannon limit



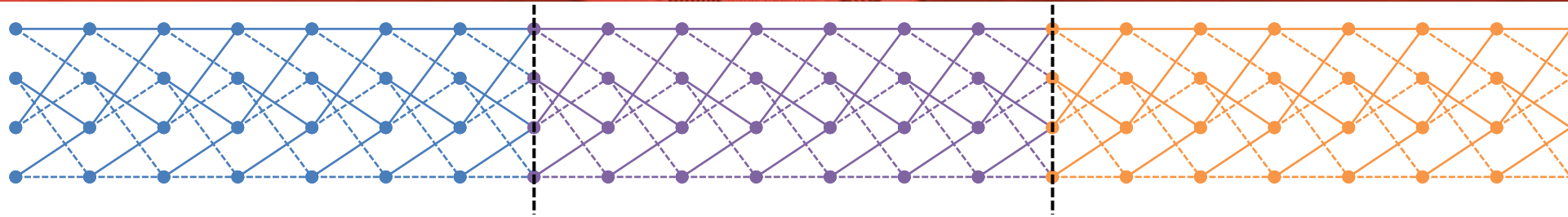
Efficient Decoding of Turbo Codes with Sliding Window BCJR

- Structurally similar to Viterbi algorithm
- Requires forward and backward iterations
 - Large memory to store intermediate results
- Decompose trellis into smaller, overlapping windows
 - Memory reduction at the cost of additional complexity
 - Forward, backward, and dummy iterations can be parallelized: less memory and higher throughput

Throughput 390 Mbps
Energy/bit ~2nJ/bit

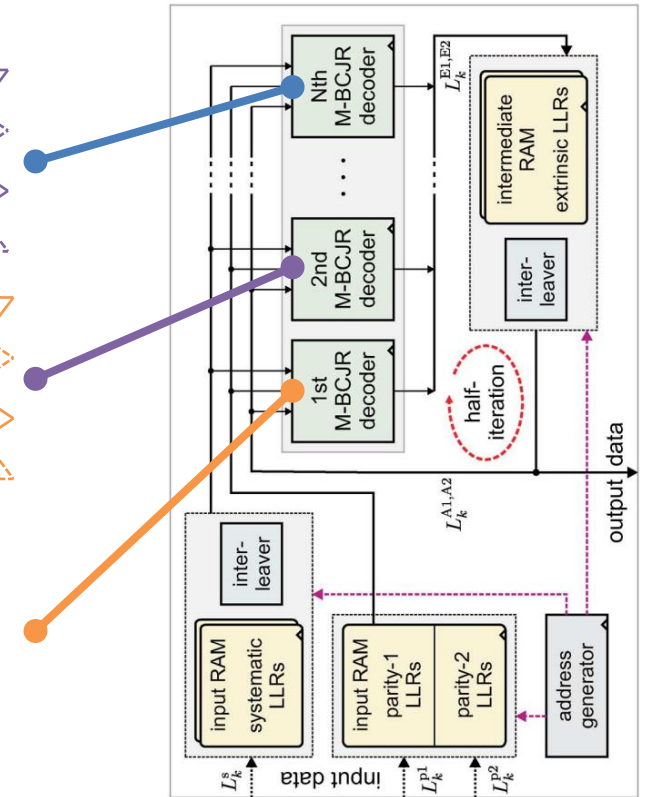
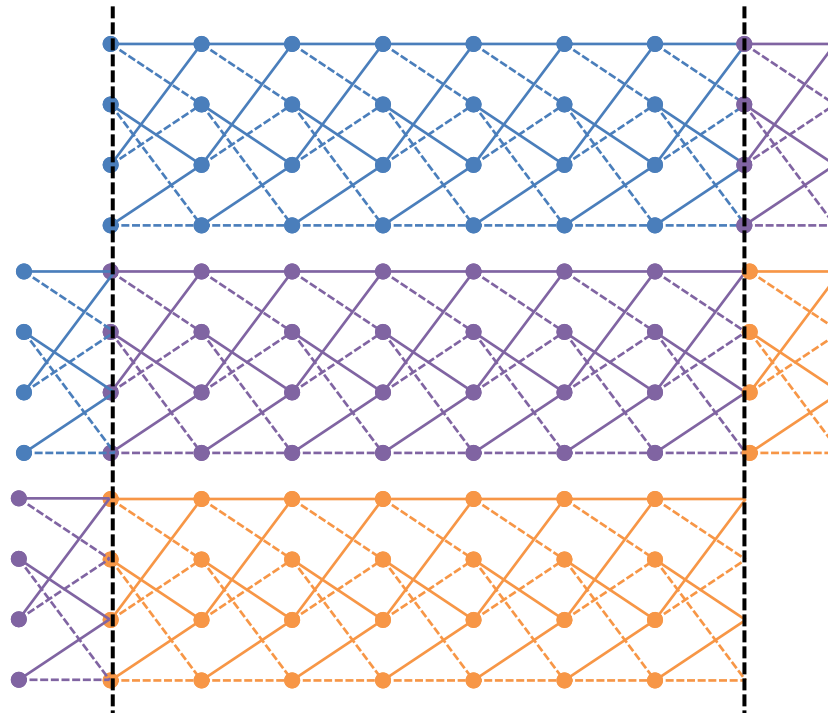


Boosting the throughput of BCJR: Parallel Sliding Window



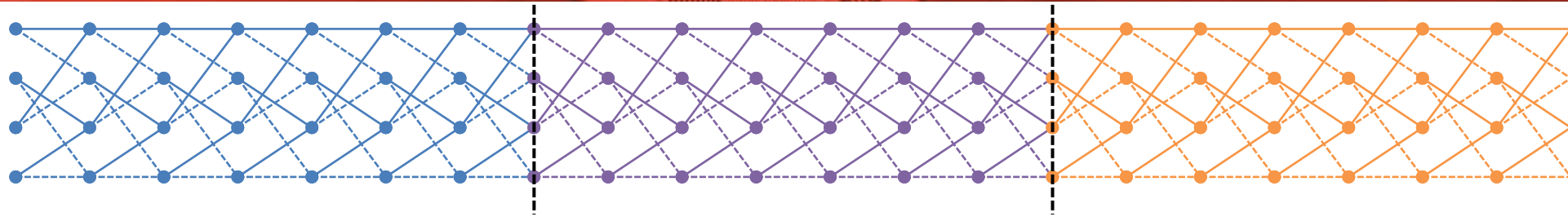
Process multiple trellis windows in parallel with parallel BCJR

- Parallel memory access
 - Even iterations (non-interleaved): no access conflicts



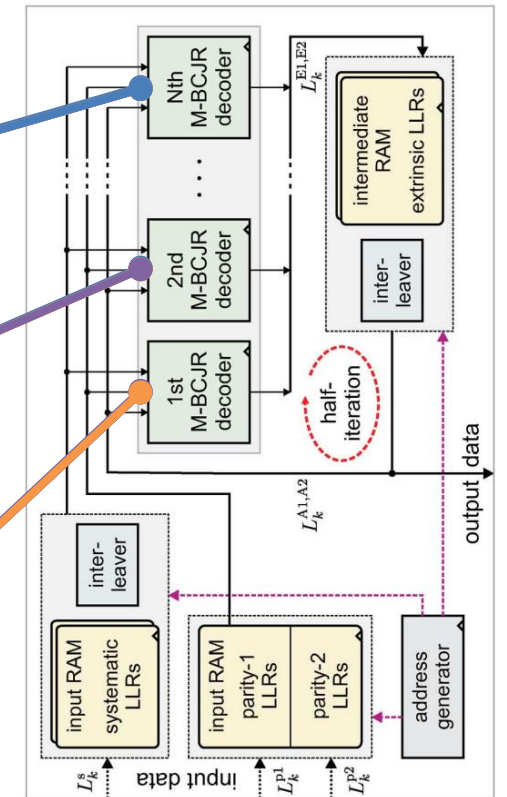
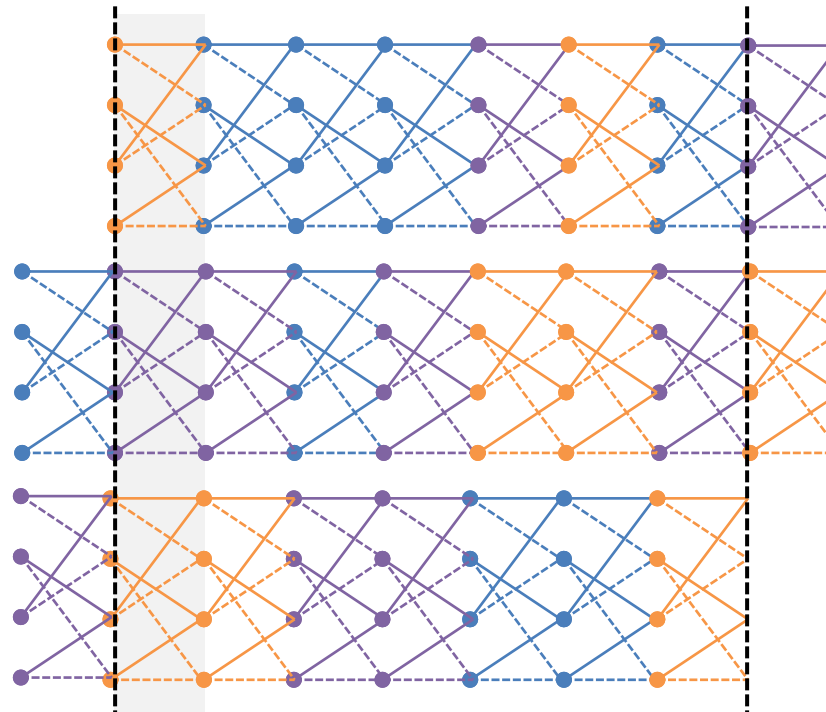
[Studer, 2011]

Boosting the throughput of BCJR: Parallel Sliding Window



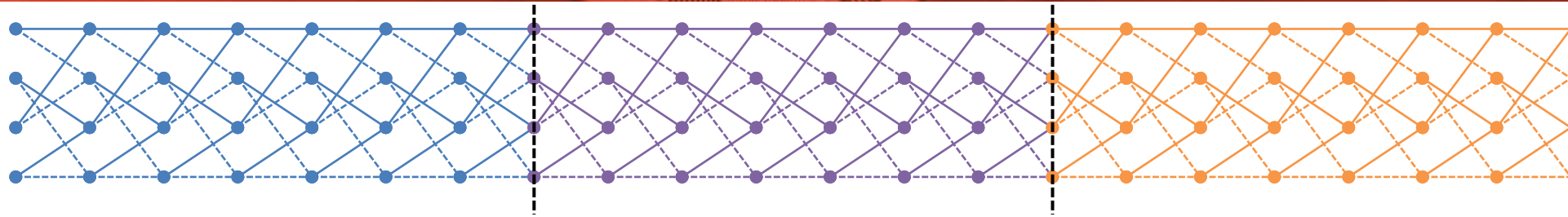
Process multiple trellis windows in parallel with parallel BCJR

- Parallel memory access
 - Non-interleaved iterations: no access conflicts
 - **Interleaved iterations: access conflicts** due to re-ordering



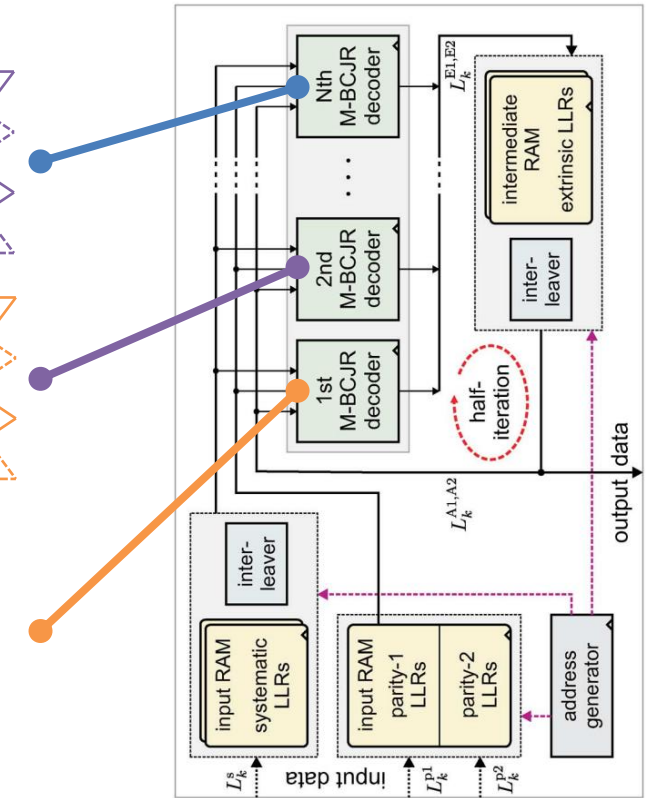
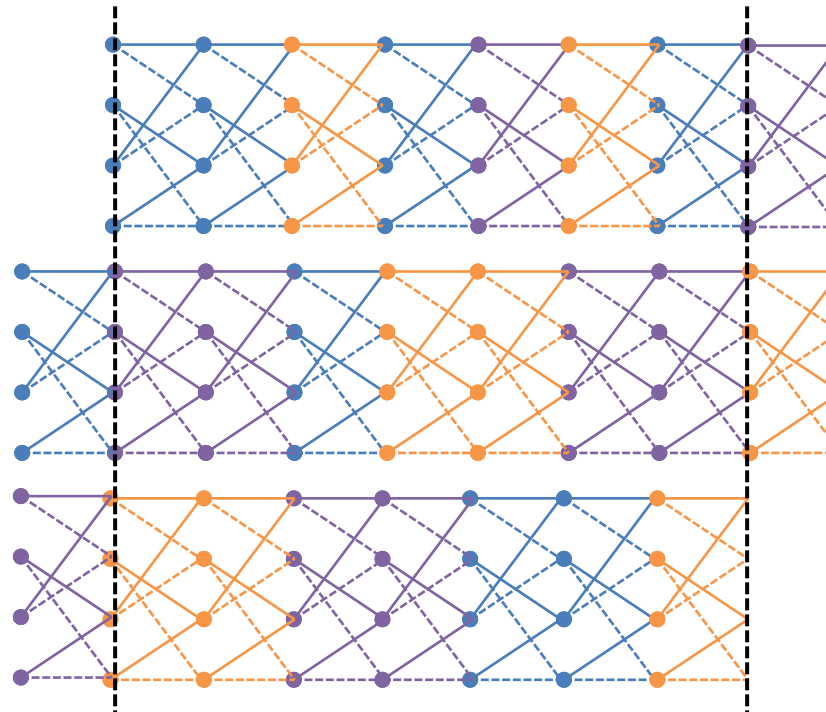
[Studer, 2011]

Boosting the throughput of BCJR: Parallel Sliding Window



Process multiple trellis windows in parallel with parallel BCJR

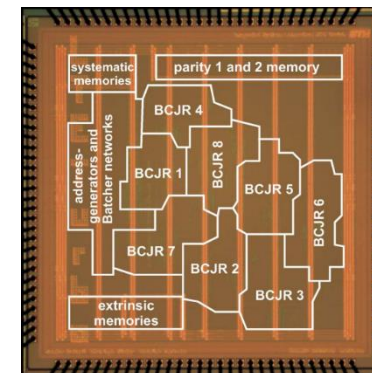
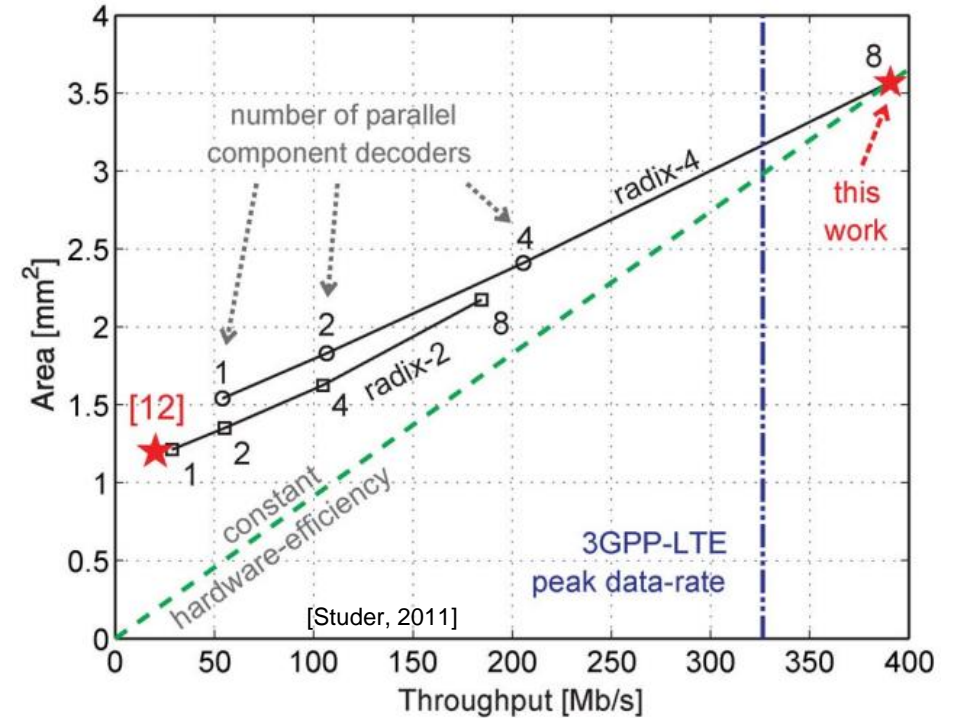
- Parallel memory access
 - Non-interleaved iterations: no access conflicts
 - **Interleaved iterations: access conflicts** due to re-ordering
- **4G-LTE** defines a hardware-optimized Quadratic-permutation-polynomial (QPP) interleaver that avoids access conflicts



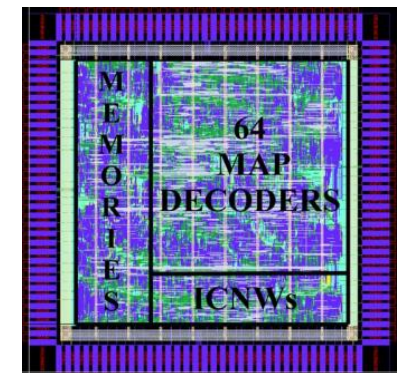
[Studer, 2011]

Boosting the throughput of BCJR: Parallel Sliding Window

- **Memory consumes a significant portion of the area**
 - More parallel functional units lead to better hardware-efficiency
- **Maximum parallelism limited by the interleaver**
 - increasing structure eventually impacts FER performance
- **Reference implementations reach to only ~1 Gbps even in advanced technology nodes**



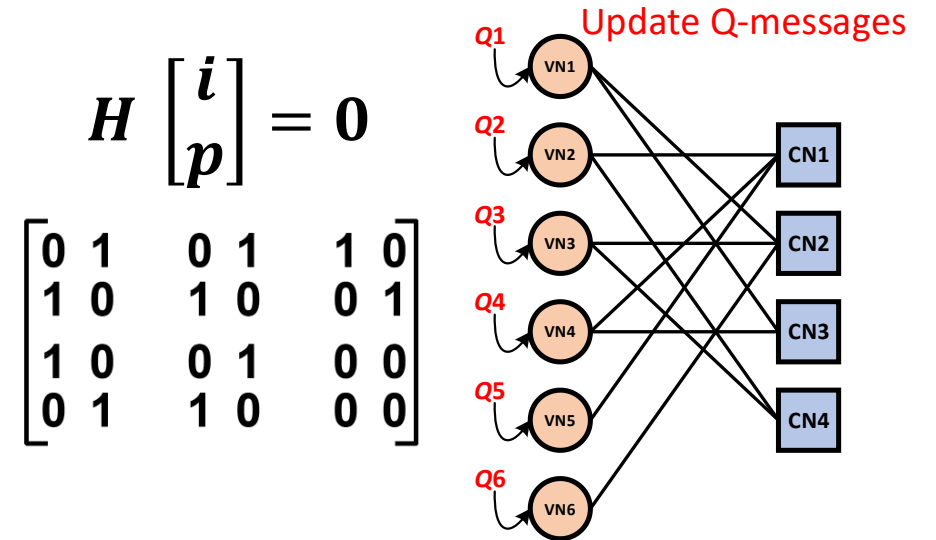
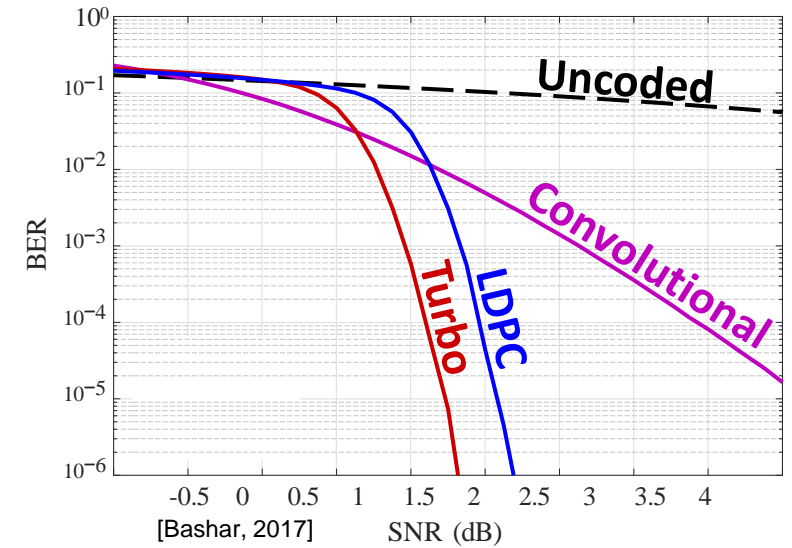
[Studer, 2011]



[R. Shrestha, 2014]

The Comeback of a Forgotten Idea: Low Density Parity Check (LDPC)

- **Discovered by Gallager already in 1963**
- Linear block code with a sparse parity check matrix
- Performance close to Shannon limit
- Belief Propagation (BP) Decoding
 - **Passing messages on a factor graph**
 - **Highly parallel**
 - **Many short component codes (CNs)**
- **Deemed to complex for implementation**
- **Come-back in the wake of VLSI parallelism (MacKey, 1996)**

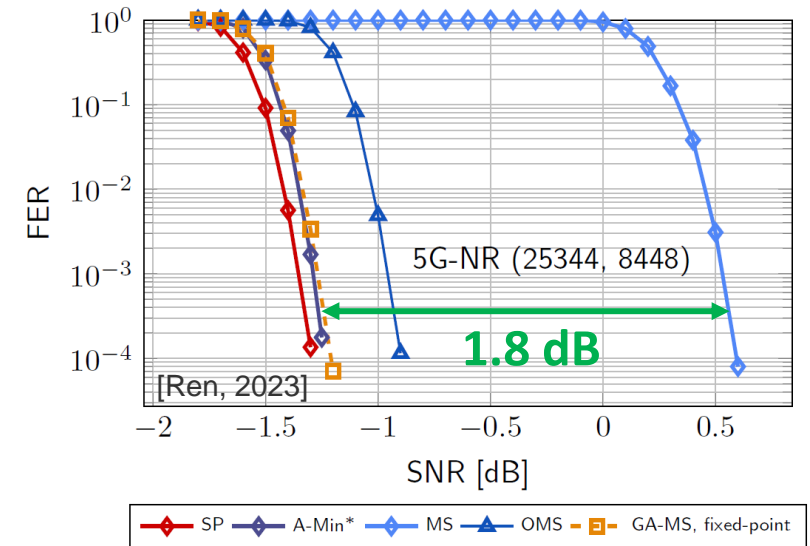


MIN-SUM and Friends: Many Approximations are Good Enough

$$R_{c,v} = \tanh^{-1} \left\{ \prod_{n' \in \mathcal{N}(m) \setminus n} \tanh(T_{v',c}/2) \right\}$$

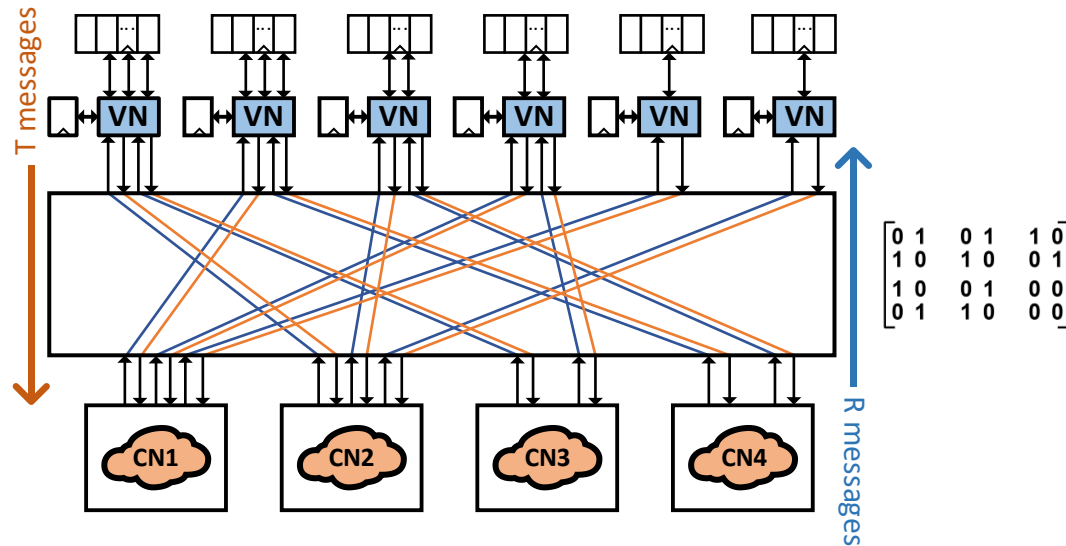
Optimal Sum-Product BP decoding has prohibitively complex message update rules

Sum Product (SP) [↙]	$R_{c,v} = \prod_{v' \in \mathcal{M}(c) \setminus v} \text{sgn}(T_{v',c}) \cdot \boxplus_{v' \in \mathcal{M}(c) \setminus v} T_{v',c} $ [↙]	F. R. Kschischang, <i>et al</i> , <i>TIT</i> , 2003 [↙]
Min-Sum (MS) [↙]	$R_{c,v} = \prod_{v' \in \mathcal{M}(c) \setminus v} \text{sgn}(T_{v',c}) \cdot \min_{v' \in \mathcal{M}(c) \setminus v} T_{v',c} $ [↙]	N. Wiberg, 1996 [↙]
Approximate Min (A-Min*) [↙]	$R_{c,v} = \begin{cases} \prod_{v' \in \mathcal{M}(c) \setminus v} \text{sgn}(T_{v',c}) \cdot \boxplus_{v' \in \mathcal{M}(c) \setminus v} T_{v',c} , & \text{if } T_{v,c} \text{ is the minimum,} \\ \prod_{v' \in \mathcal{M}(c) \setminus v} \text{sgn}(T_{v',c}) \cdot \boxplus_{v \in \mathcal{M}(c)} T_{v,c} , & \text{otherwise.} \end{cases}$ [↙]	C. Jones, <i>et al</i> , 2003 [↙]
Normalized MS (NMS) [↙]	$R_{c,v} = \alpha \cdot \prod_{v' \in \mathcal{M}(c) \setminus v} \text{sgn}(T_{v',c}) \cdot \min_{v' \in \mathcal{M}(c) \setminus v} T_{v',c} $ [↙]	J. Chen, <i>et al</i> , <i>TCOM</i> , 2005 [↙]
Offset MS (OMS) [↙]	$R_{c,v} = \prod_{v' \in \mathcal{M}(c) \setminus v} \text{sgn}(T_{v',c}) \cdot \max \left(\min_{v' \in \mathcal{M}(c) \setminus v} T_{v',c} - \beta, 0 \right)$ [↙]	J. Chen, <i>et al</i> , <i>TCOM</i> , 2005 [↙]
Self Correction MS [↙]	$T_{v,c} = 0$, if $\text{sgn}(T_{v,c}^{(i)}) \neq \text{sgn}(T_{v,c}^{(i-1)})$ [↙]	V. Savin, <i>et al</i> , <i>ISIT</i> , 2008 [↙]
Adjusted MS (A-MS) [↙]	$R_{c,v} = \begin{cases} \prod_{v' \in \mathcal{M}(c) \setminus v} \text{sgn}(T_{v',c}) \cdot \text{LUT}_{v' \in \mathcal{M}(c) \setminus v} T_{v',c} , & \text{if } T_{v,c} \text{ is the minimum,} \\ \prod_{v' \in \mathcal{M}(c) \setminus v} \text{sgn}(T_{v',c}) \cdot \text{LUT}_{v \in \mathcal{M}(c)} T_{v,c} , & \text{otherwise.} \end{cases}$ [↙]	T. Richardson, <i>et al</i> , 2018 [↙]
Adaptive MS [↙]	$R_{c,v} = \begin{cases} \prod_{v' \in \mathcal{M}(c) \setminus v} \text{sgn}(T_{v',c}) \cdot \max \left(\min_{v' \in \mathcal{M}(c) \setminus v} T_{v',c} - \beta, 0 \right), & \text{if } c \leq 4 \\ \prod_{v' \in \mathcal{M}(c) \setminus v} \text{sgn}(T_{v',c}) \cdot \min_{v' \in \mathcal{M}(c) \setminus v} T_{v',c} , & \text{otherwise} \end{cases}$ [↙]	K. Trung, <i>et al</i> , <i>ISCCAS</i> , 2019 [↙]



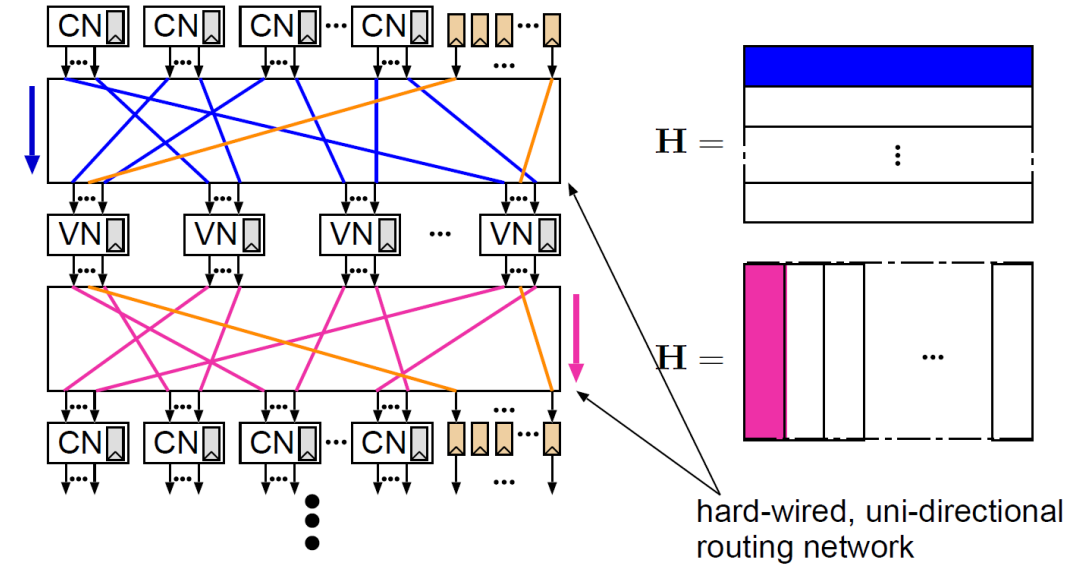
5G-NR LDPC Decoder	MS	GA-MS
T/P [Gbps]	29.40	28.25 (-3.9%)
Area [mm ²]	1.24	1.33 (+7.2%)

Scaling (More Transistors) Enable Extremely High Throughput



Isomorphic architecture

Direct mapping of Tanner graph onto silicon



Pipelined isomorphic architectures

Parallelize iterations: one cycle per code word

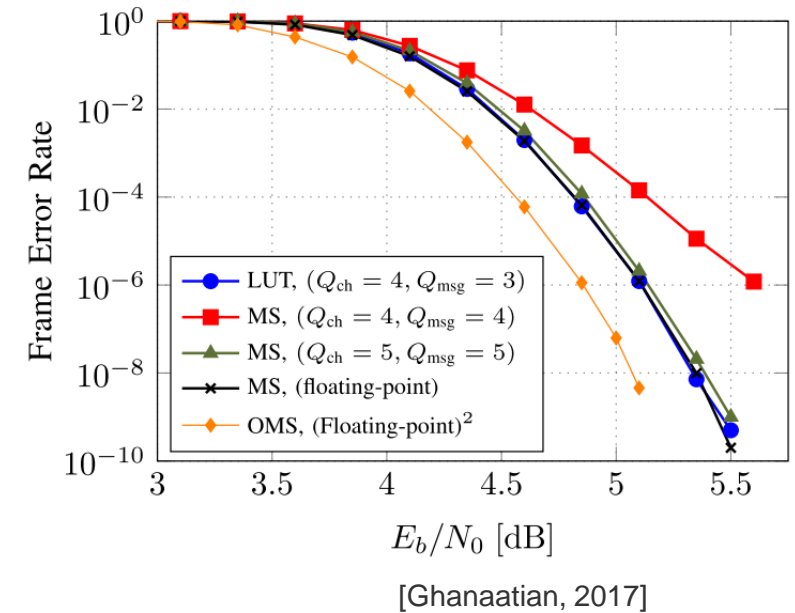
Bottleneck: Message exchange between VNs and CNs

Routing of messages between layers leads to extreme routing congestion

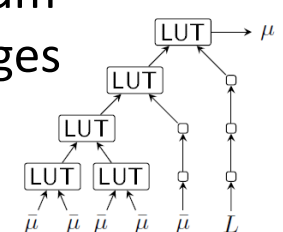
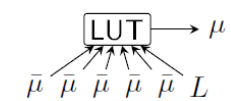
Only feasible for codes with low row/column degrees (poor performance)

Finite Alphabet Decoding to the Rescue

- **Quantization of messages in fixed-point implementations has large impact on**
 - Routing, storage, logic **complexity**
 - Decoder **error rate performance** and error floor
- Information theoretic **analysis suggests a non-uniform quantization**
 - Combining with conventional arithmetic is inefficient
- Solution: CNs and VNs based on lookup-tables
 - Information bottleneck analysis provides optimum mapping from input-messages to output messages
 - LUTs optimized for multi-stage as Boolean logic



$$\sum_{\mathbf{x}: x_0=\dots=x_{d_v-1}=x} p_{L|x}(L|x_0) \prod_{j=1}^{d_v-1} p_{m|x}^{(i)}(\bar{\mu}_j|x_j)$$



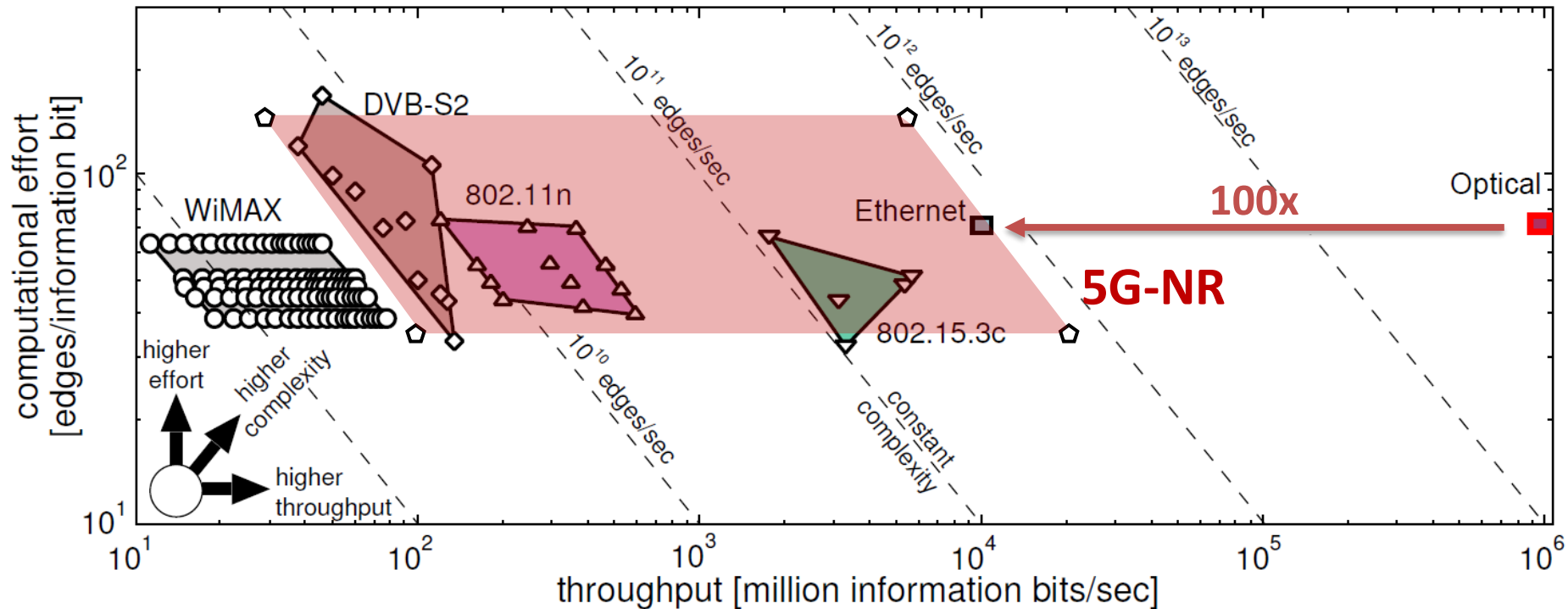
Area -30%
Throughput 2.17x
Energy/bit -50%



Throughput >500 Gbps

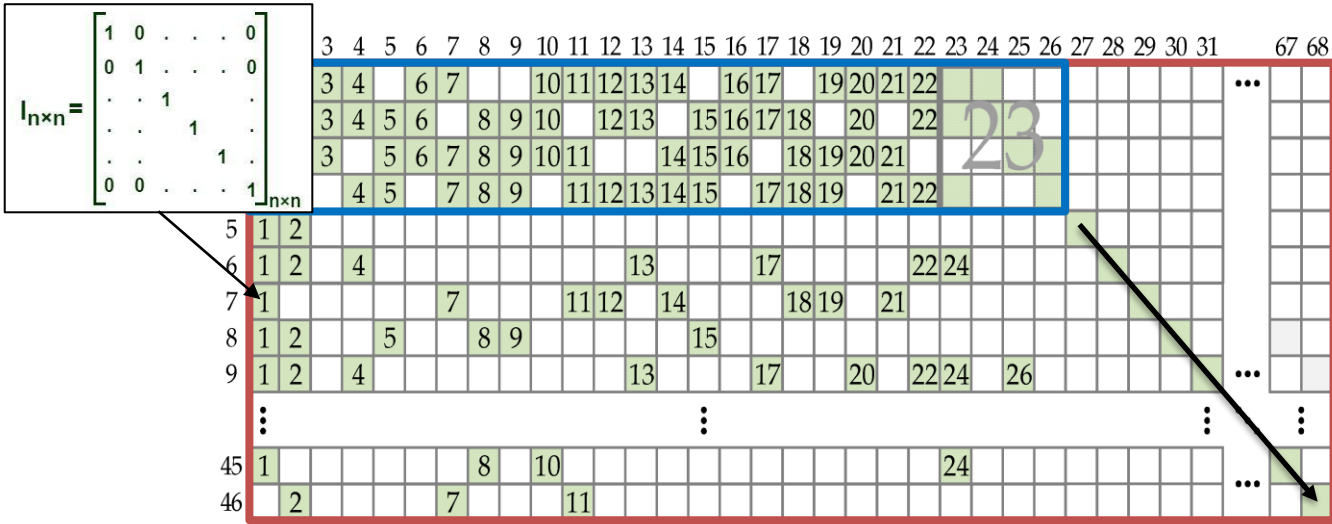
High Throughput for Optical Communications

Computational effort per information bit and required throughput determines implementation requirements



Wireless system computational requirements (e.g., 5G) are significantly lower, but require flexibility

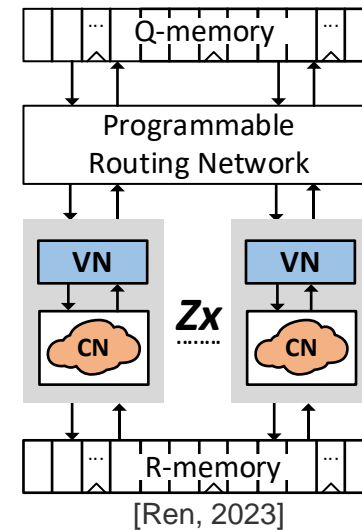
Flexible QC-LDPC Codes and Decoder Architectures for 5G



5G LDPC matrix constructed for flexibility and efficiency

Quasi cyclic code: compact base graph that expands into a larger matrix depending on block size

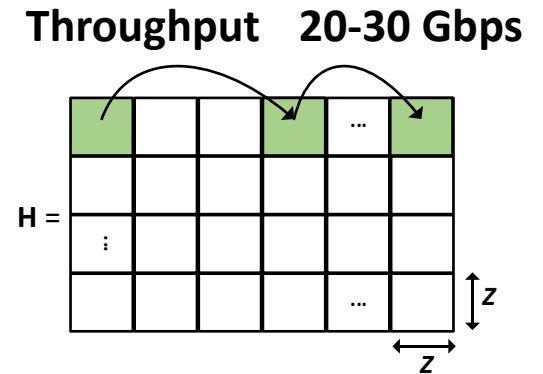
Raptor-like code: good channels (high rate) require less decoding effort than bad channels



Decoder architecture optimized for 5G QC-LDPC codes

Layered decoding: row-wise processing of the base graph (good convergence and Zx parallelism)

Block-parallel: decomposes CN and VN calculation into multiple cycles for HW efficiency



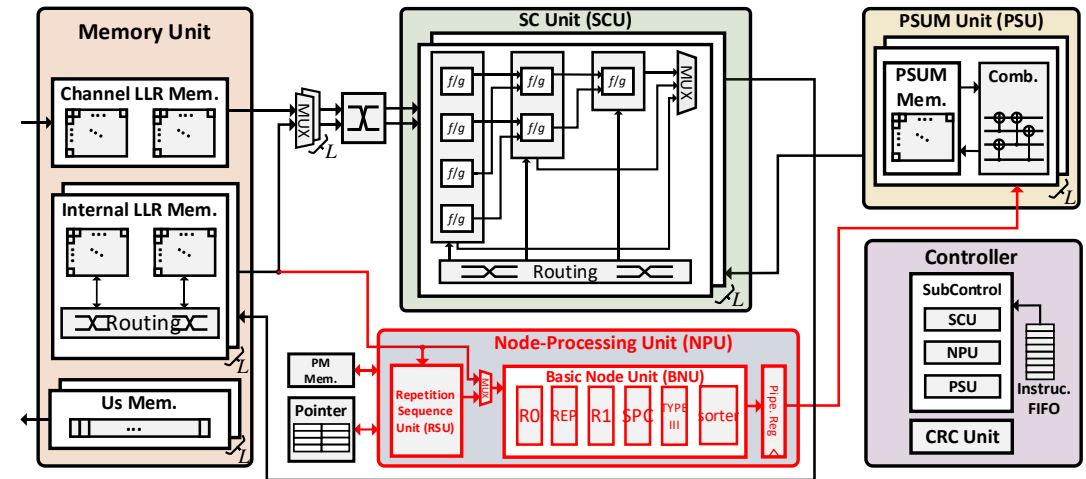
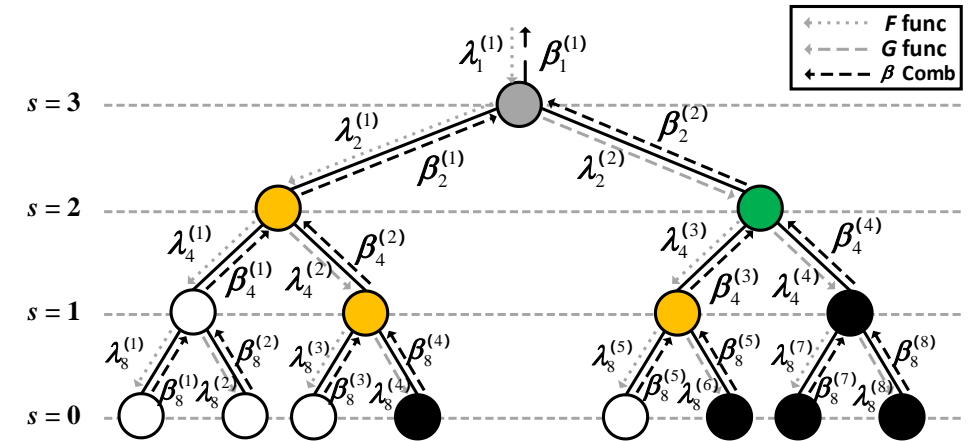
5G-NR Polar Codes for Short Block Lengths and URLLC

Polar codes: discovered by E. Arıkan

- Systematic construction with great flexibility for code length and code rate
- **Good performance for short block lengths** when concatenated with a CRC (for list decoding)
- Used for 5G-NR control channel and considered for URLLC and 6G

Decoding: successive cancellation (SC) algorithm and its list variants for better FER performance

- **Sequential process** following a tree-traversal procedure **decodes bit-by-bit**

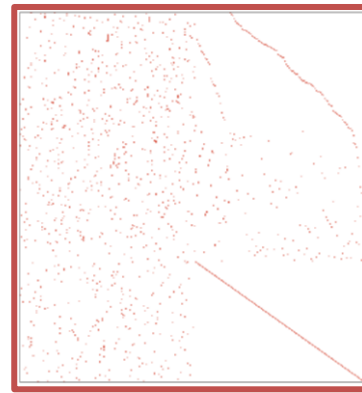
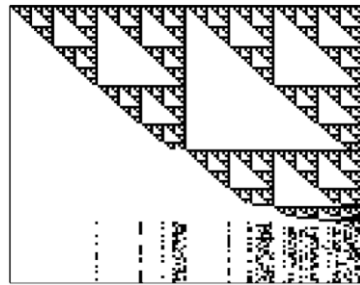
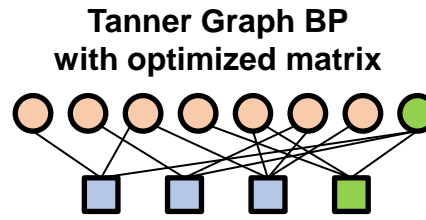
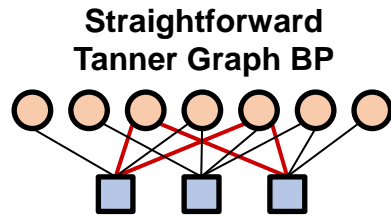
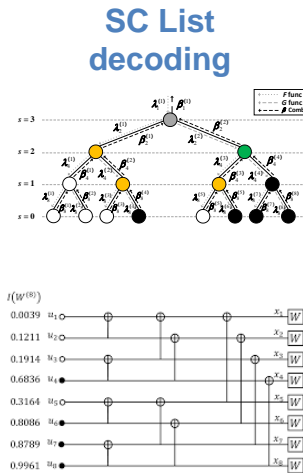


[Ren, 2022]

Throughput 3-6 Gbps

BP Decoding for Everyone

Can we decode Polar codes and other short codes with BP decoding?



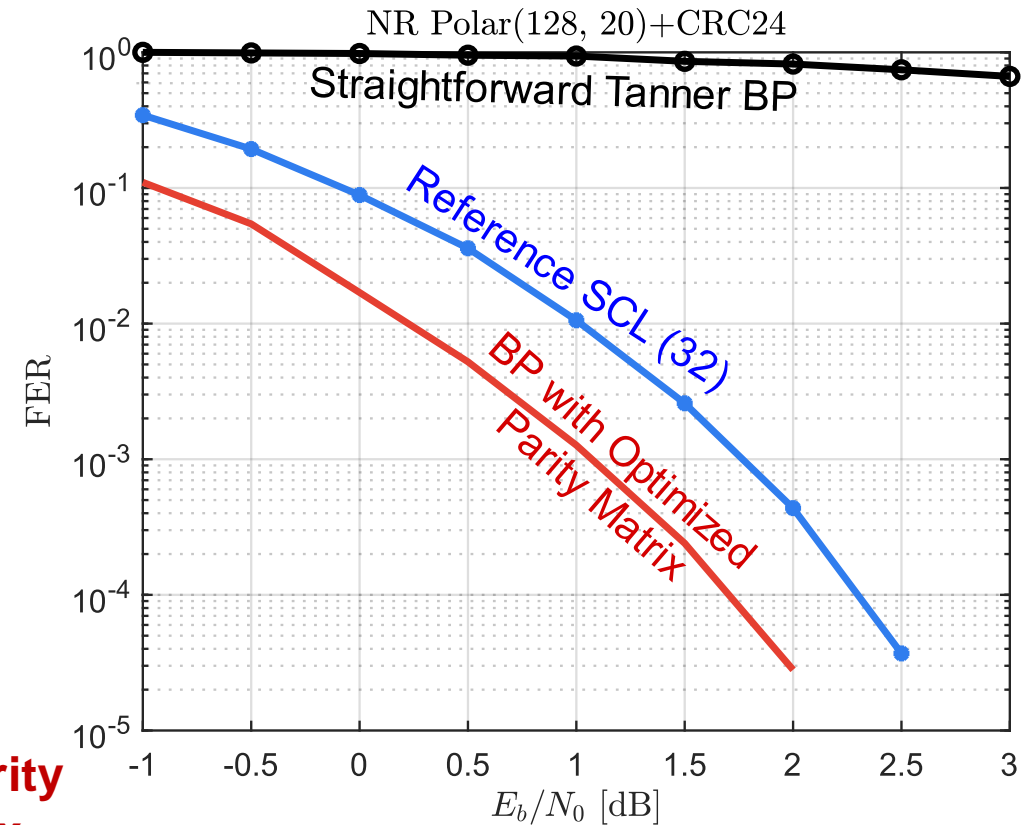
Straightforward derivation of a parity check matrix



BP decoding on a Tanner Graph as for LDPC with high parallelism

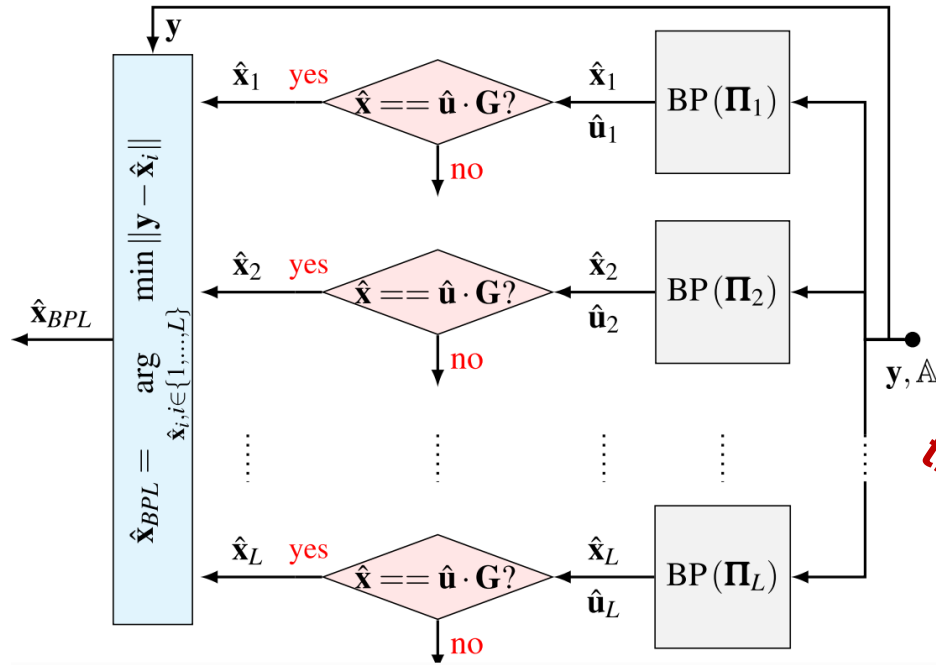
Optimized parity check matrix

Slightly larger, but more sparse with fewer cycles (better BP performance)



Close-to-Optimal Performance with Simple, Suboptimal Decoders

Key insight: a given code can be represented by many different parity check matrices

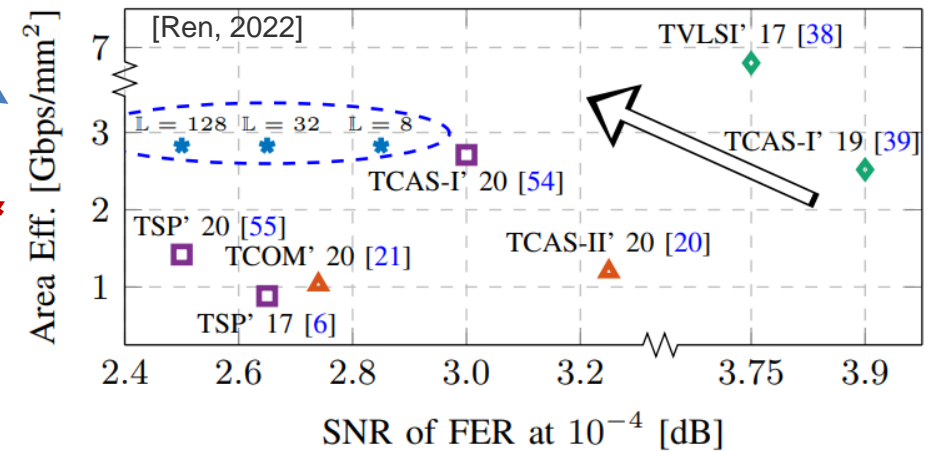
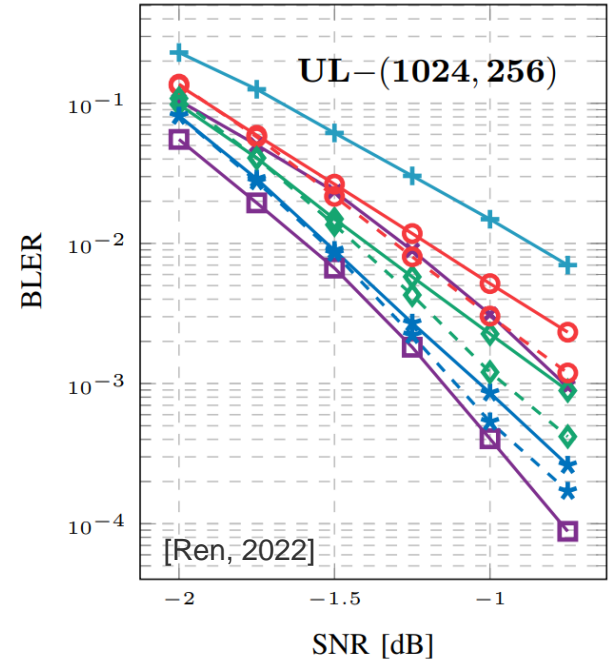


Example: BP List Decoding [Elkelesh, 2018]

Automorphism Ensemble Decoding (AED):
multiple sub-optimal decoders try to decode with different parity check matrices

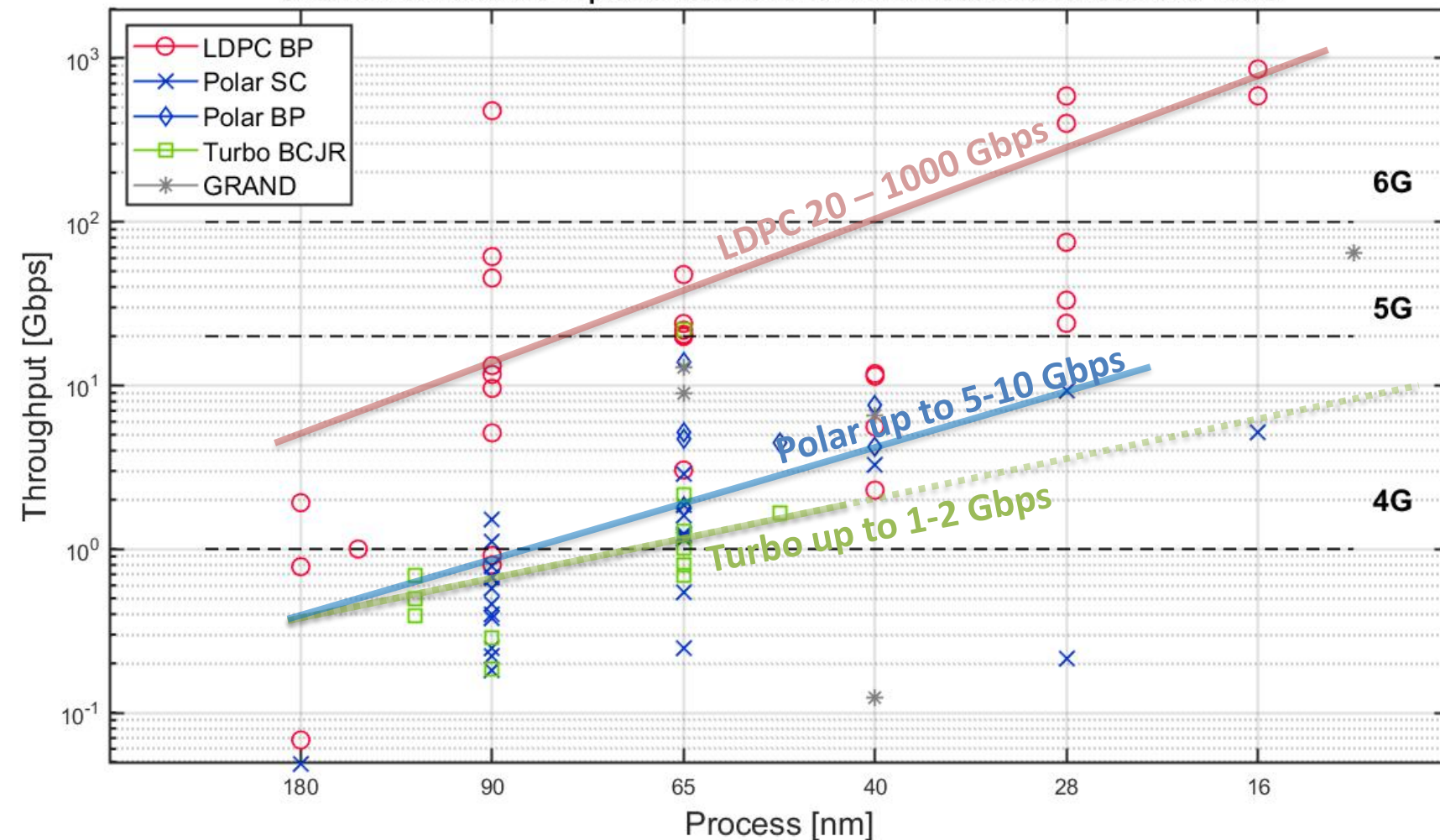
More attempts, better performance

Many sub-optimal decoders more efficient than one optimal decoder



Conclusions

Overview of Hardware Implementations for Channel Decoders in Recent 20 Years



- Technology continues to scale, but improvements are mostly due to density
- Algorithm choices and optimizations are key to efficient implementation
- Parallel processing is essential to exploit the remains of Moore's law
 - Sequential decoders fall short from scaling
 - BP decoding scales well with technology
- Memory becomes an increasingly serious issue

Acknowledgements



Yuqing
Ren



Yifei
Shen



Andreas
Toftegaard
Kristensen



Reza
Ghanaatian



Alexios
Balatsoukas-
Stimming
