

A 354 Mb/s 0.37 mm² 151 mW 32-User 256-QAM Near-MAP Soft-Input Soft-Output Massive MU-MIMO Data Detector in 28nm CMOS

Charles Jeon, Oscar Castañeda, and Christoph Studer

Abstract—This paper presents a novel data detector ASIC for massive multiuser multiple-input multiple-output (MU-MIMO) wireless systems. The ASIC implements a modified version of the large-MIMO approximate message passing algorithm (LAMA), which achieves near-optimal error-rate performance (i) under realistic channel conditions and (ii) for systems with as many users as base-station (BS) antennas. The hardware architecture supports 32 users transmitting up to 256-QAM simultaneously and in the same frequency band, and provides soft-input soft-output capabilities for iterative detection and decoding. The fabricated 28nm CMOS ASIC occupies 0.37 mm², achieves a throughput of 354 Mb/s, consumes 151 mW, and improves the SNR by more than 11 dB compared to existing data detectors in systems with 32 BS antennas and 32 users for realistic wireless channels. In addition, the ASIC achieves 4× higher throughput per area than a recently proposed message-passing detector.

I. INTRODUCTION

Massive MU-MIMO enables higher per-cell spectral efficiency compared to conventional, small-scale MIMO. This improvement, however, comes at a significant increase in baseband processing complexity [1]. In particular, data detection at the base-station (BS) in the massive MU-MIMO uplink is among the most critical tasks in terms of power consumption and throughput [2]. To exacerbate the situation, the complexity of optimal, maximum a-posteriori (MAP), data detection grows exponentially in the number of user equipment (UE) antennas [3], which prevents its implementation in practice.

To enable high-throughput massive MU-MIMO data detection, a variety of low-complexity algorithms (see, e.g., [1], [4]) and application-specific integrated circuits (ASICs) [5]–[7] have been proposed. These algorithms and ASIC designs either rely on idealistic channel-hardening assumptions [5], [6] or deploy approximations [1], [4] to reduce complexity. Unfortunately, both of these simplifications result in high error rates (i) under realistic propagation conditions, such as correlation and per-user path loss, and (ii) in systems where the number of UEs is equal to the number of BS antennas. As a consequence, achieving near-optimal performance in realistic systems necessitates novel data detection algorithms that can be implemented efficiently.

CJ, OC, and CS are with the School of ECE, Cornell University, Ithaca, NY; e-mail: studer@cornell.edu; web: <http://vip.ece.cornell.edu>. The work was supported by Xilinx Inc. and by the US NSF under grants ECCS-1408006, CCF-1535897, CCF-1652065, CNS-1717559, and ECCS-1824379.

The authors thank A. Maleki for discussions on approximate message passing and F. K. Gürkaynak for his assistance during ASIC testing.

Contributions: We propose the first data detector ASIC that achieves near-MAP performance for 32 UEs under realistic propagation conditions. Furthermore, the ASIC provides soft-input soft-output (SISO) capabilities for iterative detection and decoding. The algorithm builds upon the large-MIMO approximate message passing (LAMA) algorithm [3], which achieves MAP-optimal error-rate performance for Rayleigh fading channels and in the large-antenna limit, assuming that the UE-to-BS antenna ratio is less than a threshold that depends on the constellation. In contrast to linear data detectors, LAMA exploits information on the constellation to improve performance; for QPSK, for example, LAMA achieves optimal performance in the large-antenna limit and for systems where the number of UE and BS antennas are identical. Since practical systems are finite-dimensional and real-world channels exhibit correlation, we include algorithm-level optimizations to support realistic channels with LAMA. To achieve high throughput at low area, our ASIC uses coarse-grained pipeline interleaving, processing two detection problems within the same architecture. The fabricated 28nm CMOS ASIC outperforms existing designs under realistic channel conditions and for systems in which the number of UEs is comparable to the number of BS antennas.

II. MASSIVE MU-MIMO DATA DETECTION

We consider the uplink of a coded massive MU-MIMO system with U single-antenna UEs and B BS antennas. The information bit vectors \mathbf{b} of the U UEs are encoded on a per-UE basis (e.g., using a convolutional code) and the resulting coded bit-stream vectors \mathbf{x} are mapped (using Gray labeling) to a sequence of transmit vectors $\mathbf{s} \in \mathcal{O}^U$, where \mathcal{O} corresponds to the constellation of size 2^Q . Each transmit vector \mathbf{s} is associated with UQ binary values $x_{u,q} \in \{0, 1\}$, $u = 1, \dots, U$, $q = 1, \dots, Q$, corresponding to the q th bit of the u th entry (spatial stream) of \mathbf{s} . We assume $\mathbb{E}_s[\mathbf{s}\mathbf{s}^H] = E_s \mathbf{I}_U$, where E_s is the symbol variance. The baseband input-output relation of the MU-MIMO channel is modeled as $\mathbf{y} = \mathbf{H}\mathbf{s} + \mathbf{n}$, where $\mathbf{H} \in \mathbb{C}^{B \times U}$ is the MIMO channel matrix, $\mathbf{y} \in \mathbb{C}^B$ is the received vector at the BS, and \mathbf{n} is B -dimensional i.i.d. zero-mean complex Gaussian distributed noise with variance N_0 per entry. We assume that \mathbf{H} , N_0 , and E_s are known at the BS.

A. Iterative MIMO Decoding

Iterative detection and decoding in MIMO systems achieves near-optimal spectral efficiency in MIMO wireless systems [8].

Algorithm 1 Large MIMO AMP (LAMA) Algorithm

- 1: **inputs:** \mathbf{H} , \mathbf{y} , N_0 , and $\Lambda_{u,q}^{\text{prior}}, \forall u, q$
- 2: **preprocessing:** $\tilde{\mathbf{G}} = \mathbf{I}_U - \text{diag}(\mathbf{G})^{-1}\mathbf{G}$ with $\mathbf{G} = \mathbf{H}^H\mathbf{H}$, $\tilde{\mathbf{y}}^{\text{MF}} = \text{diag}(\mathbf{G})^{-1}\mathbf{H}^H\mathbf{y}$, and $g_u = G_{uu}/U, u = 1, \dots, U$
- 3: **initialize:** $\mathbf{z}^1 = \hat{\mathbf{s}}^1 = \mathbf{0}_{U \times 1}$, and $\rho^1 = 0$
- 4: **for** $t = 1, 2, \dots, t_{\text{max}}$ **do**
- 5: **mean and variance estimation:**

$$\hat{\mathbf{s}}^{t+1} = \mathbf{F}(\mathbf{z}^t, \rho^t, \mathbf{g}, \Lambda^{\text{prior}}) \quad (\text{mean update})$$

$$\boldsymbol{\tau}^{t+1} = \mathbf{G}(\mathbf{z}^t, \rho^t, \mathbf{g}, \Lambda^{\text{prior}}) \quad (\text{variance update})$$

$$\hat{\boldsymbol{\tau}}^{t+1} = \frac{1}{B} \mathbf{g}^T \boldsymbol{\tau}^{t+1}$$

$$\boldsymbol{\alpha}^t = \mathbf{z}^t - \hat{\mathbf{s}}^t \quad (\text{Onsager term})$$

$$\mathbf{b}^t = \rho^t \hat{\boldsymbol{\tau}}^{t+1}$$
- 6: **interference cancellation:**

$$\mathbf{z}^{t+1} = \tilde{\mathbf{y}}^{\text{MF}} + \tilde{\mathbf{G}}\hat{\mathbf{s}}^{t+1} + \mathbf{b}^t \boldsymbol{\alpha}^t \quad (\text{interference cancellation})$$

$$\rho^{t+1} = \left(\frac{1}{B} N_0 + \hat{\boldsymbol{\tau}}^{t+1}\right)^{-1} \quad (\text{post-equalization SINR update})$$
- 7: **end for**
- 8: **output:** extrinsic LLR values $\Lambda_{u,q}^{\text{d}}, \forall u = 1, \dots, U, q = 1, \dots, Q$

Reliability information on the coded bits, often expressed as log-likelihood ratios (LLRs), is iteratively exchanged between the MIMO data detector and the channel decoder. In each iteration, a soft-input soft-output (SISO)-capable MIMO data detector computes extrinsic LLRs for the coded bits $x_{u,q}$ as

$$\Lambda_{u,q}^{\text{d}} = \log \left(\frac{\mathbb{P}[x_{u,q} = 1 | \mathbf{y}]}{\mathbb{P}[x_{u,q} = 0 | \mathbf{y}]} \right) - \Lambda_{u,q}^{\text{prior}},$$

using the received vector \mathbf{y} and a-priori LLRs $\Lambda_{u,q}^{\text{prior}}, u = 1, \dots, U, q = 1, \dots, Q$, obtained from the channel decoder. The extrinsic LLRs $\Lambda_{u,q}^{\text{d}}$, which represent reliability estimates for each coded bit $x_{u,q}$, are then passed to the channel decoder, which computes *new* a-priori LLRs $\Lambda_{u,q}^{\text{prior}}, \forall u, q$, that are used by the MIMO data detector in the next iteration. After a small number of iterations I , the channel decoder generates final decisions $\hat{\mathbf{b}}$ for the information bit vector \mathbf{b} .

B. Hardware Friendly LAMA Algorithm

LAMA is an efficient data detection algorithm based on approximate message passing (AMP), that is provably optimal (in terms of error-rate performance) in the large-system limit (i.e., fix $\beta = B/U$ and $B \rightarrow \infty$) with i.i.d. Rayleigh fading channels [3]. In each of its t_{max} iterations, LAMA decouples the MIMO system into parallel and independent AWGN channels with equal signal-to-interference-plus-noise ratio (SINR). As a result, LAMA optimally denoises the parallel AWGN channels in every iteration, which successively increases the post-equalization SINR and improves the error-rate performance.

To deal with realistic channel conditions (such as correlation and per-UE path loss), we apply algorithm-level modifications to the original LAMA algorithm in [3]. First, we transform LAMA so that it operates on the $U \times U$ dimensional Gram matrix $\mathbf{G} = \mathbf{H}^H\mathbf{H}$, instead of the $B \times U$ channel matrix \mathbf{H} , which reduces the per-iteration complexity. Second, we deploy message damping techniques [9] to reduce the performance loss of LAMA in finite-dimensional systems that exhibit correlation and large-scale UE fading. Specifically, we damp the updates of $\hat{\boldsymbol{\tau}}^t$ and ρ^t by a factor $\theta \in (0, 1]$, i.e., we use $\hat{\boldsymbol{\tau}}_d^t$ instead of $\hat{\boldsymbol{\tau}}^t$ in line 6 of Algorithm 1, and $\hat{\boldsymbol{\tau}}_d^t = \theta \hat{\boldsymbol{\tau}}^t + (1 - \theta) \hat{\boldsymbol{\tau}}_d^{t-1}$. Third, we include support for iterative detection and decoding. The implemented LAMA algorithm is summarized in Algorithm 1

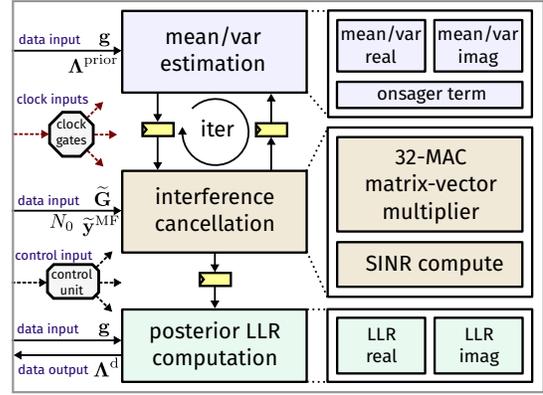


Fig. 1. Top-level architecture of the LAMA data detector. Data detection is carried out in a pipeline-interleaved manner, which iteratively processes two independent detection problems in the mean and variance (MV) estimation unit and the interference cancellation (IC) unit.

(message damping details are excluded). The functions \mathbf{F} and \mathbf{G} correspond to the posterior mean and variance applied element-wise, i.e., $\mathbf{F}(z, \rho, \Lambda^{\text{prior}}) = \mathbb{E}_S(S|z = S + \rho^{-1/2}N)$, $N \sim \mathcal{CN}(0, 1)$ and $p(S)$ can be derived from the a-priori LLRs Λ^{prior} ; \mathbf{G} can be derived similarly—see [3] for the details.

III. VLSI ARCHITECTURE

Fig. 1 depicts the top-level architecture of the LAMA data detector. LAMA performs two main tasks per iteration: The first task estimates the mean and variance (MV) of the data transmitted by each UE; the second task cancels interference (IC) among the UEs—both of these tasks are detailed below. To maximize throughput, two independent detection problems are processed simultaneously in a pipeline-interleaved manner, i.e., one problem per task. The two main processing units, namely MV and IC, perform the assigned computations in T_s clock cycles, and the results of both units are exchanged for further processing in the subsequent iteration. In the last t_{max} iteration, the outputs from the IC unit are sent to the LLR computation unit, which takes T_{LLR} clock cycles. Thus, LAMA delivers a new set of UQ LLR values at a sustained throughput of

$$\Theta = \frac{UQ}{t_{\text{max}}T_s + T_{\text{LLR}}} f_{\text{clk}} \quad [\text{bit/s}]. \quad (1)$$

The final design supports 32 UEs, which requires $T_s = 36$ clock cycles and $T_{\text{LLR}} = 1$ clock cycle.

A. Mean and Variance Estimation (MV) Unit

In the first task (line 5 in Algorithm 1), the MV unit receives estimates of the UE's data and the associated SINR to compute mean and variance values. As shown in Fig. 2, a straightforward MV unit would require a large number of multipliers. Although statistical independence in the real and imaginary parts of the transmitted constellation points simplifies computation from M^2 -QAM to two M -PAM constellations, mean and variance computation for 16-PAM (to support 256-QAM) still requires 16 likelihood function units and a division, resulting in high complexity. Furthermore, the intrinsic LLR values obtained from the channel decoder must be transformed from bit-domain to symbol-domain for SISO processing. To

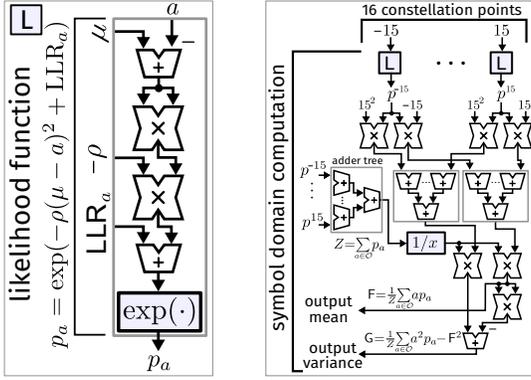


Fig. 2. Original MV unit for 256-QAM with separation into two 16-PAM units. Exact mean and variance computation entails high complexity and requires high arithmetic precision.

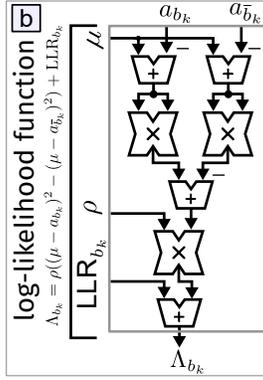


Fig. 3. Proposed low-complexity computation of message mean and variance. We transform all computations into the bit-domain and use the max-log approximation.

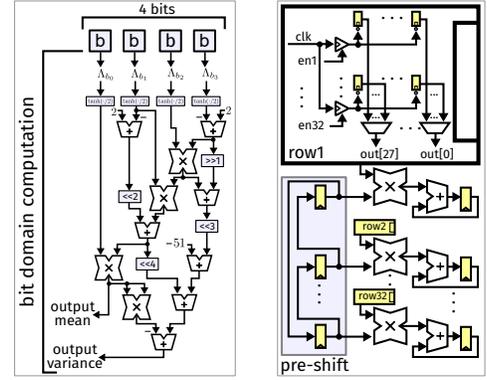


Fig. 4. Fan-out reduction of matrix-vector multiplication using Cannon's algorithm [10].

reduce complexity, existing ASICs [5]–[7] use hard-symbol clipping or linear approximations, and do not provide support for SISO processing. However, accurate message mean and variance computation is key to support realistic channels and systems with a comparable number of UEs and BS antennas.

To accurately compute the message mean and variance at low complexity, we (i) compute all quantities in the bit-domain, (ii) exploit Gray-mapping symmetries, and (iii) use the max-log approximation. The conversion into the bit-domain and the max-log approximation only requires 4 log-likelihood functions for 16-PAM, instead of 16 functions in the symbol domain. In addition, we avoid the need of a division per UE by using a LUT-based $\tanh(\cdot)$ function as in [8] with 7 input bits.

The resulting architecture, depicted in Fig. 3, also avoids the need of a division per UE. Furthermore, the architecture naturally supports SISO processing for iterative MIMO decoding. Our simulations in Section IV for various antenna configurations and channel models show that our approach entails a negligible performance loss at around $4\times$ lower area.

B. Interference Cancellation (IC) Unit

In the second task (line 6 in Algorithm 1), the IC unit performs interference cancellation and updates the SINR.

1) *32-MAC matrix-vector multiplication*: Interference cancellation requires a 32×32 complex-valued matrix-vector multiplication, which we compute sequentially in 32 clock cycles using a linear array of 32 complex-valued multiply-accumulate (MAC) units in a column-by-column fashion. To minimize the critical path caused by the large fan-out of a conventional linear array of MAC units, Fig. 4 shows a simplified version of Cannon's algorithm [10], which circularly shifts the array's input vector while sequentially processing rows of the matrix over multiple clock cycles; this reduces the vector memory fan-out from 32 MAC units to one MAC unit and a register. To further reduce the critical path and simplify placement, each row of the Gram matrix is stored next to each MAC unit with standard-cell-based latch-arrays.

2) *SINR computation*: The post-equalization SINR is computed in parallel using a Newton-Raphson (NR) reciprocal

unit [8]. We first shift the input x according to $\bar{x} = 2^\alpha x$, $\alpha \in \mathbb{Z}$ so that $\bar{x} \in [0.5, 1)$, resulting in high numerical stability. Based on an initial guess obtained from a look-up table, a single NR iteration is sufficient to compute $\bar{y}_1 \simeq \bar{x}^{-1}$; the final result $y = 2^\alpha \bar{y}_1$ corresponds to an approximation of x^{-1} .

IV. IMPLEMENTATION RESULTS AND COMPARISON

Figs. 5 and 6 show the PER of our LAMA ASIC in comparison with the linear minimum-mean squared-error (MMSE) equalizer and channel hardening-exploiting message passing (CHEMP) algorithm [4]. The number of algorithm iterations are indicated after the dash; e.g., LAMA-14 represents LAMA with 14 iterations. Outer iterations over the channel decoder are shown as either none ($I = 0$; solid lines) or one ($I = 1$; dashed lines) iteration. We simulate an LTE-based massive MU-MIMO-OFDM system at $f_c = 2$ GHz with 1200 active subcarriers and per-user convolutional coding with rate R . We use two channel models: (a) Rayleigh fading and (b) WINNER II typical urban micro [11] to model a realistic propagation environment. For a typical 256×32 ($B \times U$) massive MU-MIMO scenario, LAMA achieves the same performance as linear MMSE, but avoids a matrix inversion; CHEMP suffers an error floor above 10% PER. For the challenging 32×32 system, LAMA significantly outperforms the linear MMSE detector, achieving more than 11 dB SNR improvements for the typical urban micro channel; CHEMP fails to successfully detect packets. Extensive numerical simulations have been carried out to determine the ASIC's fixed-point parameters; the implemented design achieves near-floating-point performance.

A. Implementation Results

Fig. 7 shows a micrograph of the fabricated and fully-functional 28nm CMOS ASIC with the LAMA detector core highlighted. The LAMA ASIC only occupies 0.37 mm^2 ; the rest of the chip contains unrelated designs. The clock signal was generated by a VLSI test system and directly fed into the ASIC. At nominal supply of 0.9 V at 300 K, the ASIC reaches a maximum measured clock frequency of 400 MHz at 151 mW, which results in 354 Mb/s for 32 UEs transmitting 256-QAM.

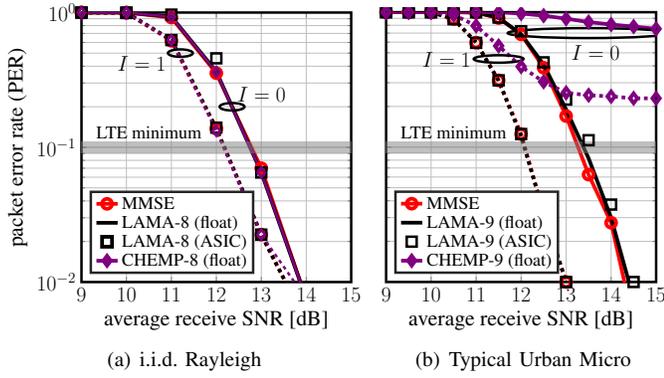


Fig. 5. 256×32 massive MU-MIMO; $R = 0.5$; 256-QAM; 9600 bits/packet.

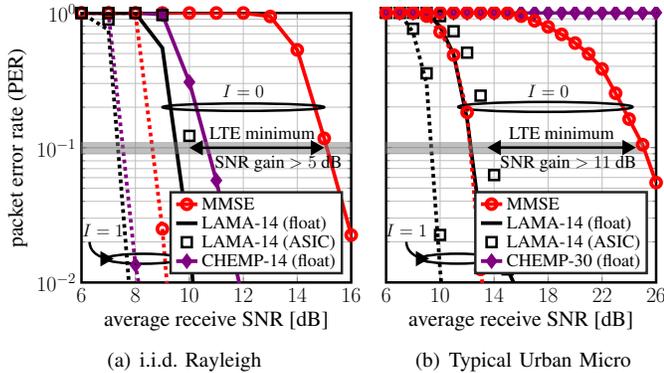


Fig. 6. 32×32 massive MU-MIMO; $R = 0.75$; QPSK; 3600 bits/packet.

Fig. 8 shows measured energy-efficiency in pJ/bit obtained via voltage-frequency scaling. By reducing the supply close to the threshold voltage, the detector achieves optimal energy-efficiency: at 0.35 V we have 123 pJ/bit (achieving 2.66 Mb/s). If maximum throughput is desired, one can increase the supply to 1.15 V and obtain 511 Mb/s (at 670 pJ/b efficiency).

Table I compares LAMA to state-of-the-art massive MU-MIMO data detectors. Our LAMA ASIC achieves more than $4\times$ improved normalized area efficiency than [7], which computes a matrix inversion. Although LAMA achieves lower area efficiency (in Gb/s/mm²) than the detectors in [5], [6], these designs suffer an error floor higher than LTE specifications under realistic channel conditions (cf. Figs 5 and 6). We note that the nominal energy efficiency is inferior to other designs due to increased arithmetic precision requirements in support of realistic channel conditions and symmetric massive MU-MIMO systems. To the best of our knowledge, the proposed LAMA ASIC is the first silicon prototype of a 32-UE massive MU-MIMO data detector that provides near-optimal error rates under realistic propagation conditions and for symmetric systems. Both of these advantages are critical to BS providers as one can support up to 32 UEs with relatively small ($B \geq 32$) BS antenna arrays under realistic channel conditions.

REFERENCES

[1] E. Larsson, O. Edfors, F. Tufvesson, and T. Marzetta, “Massive MIMO for next generation wireless systems,” *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 186–195, Feb. 2014.

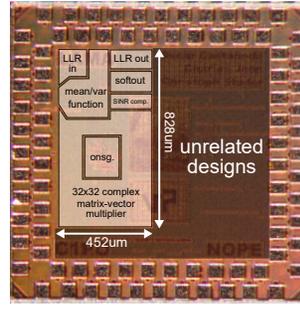


Fig. 7. Chip micrograph; LAMA data detector ASIC is highlighted.

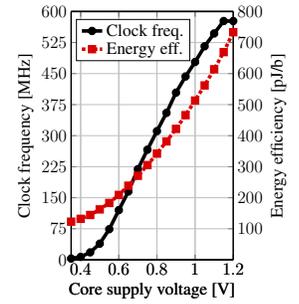


Fig. 8. Measured frequency and energy for different core voltages.

TABLE I
PERFORMANCE SUMMARY AND ASIC COMPARISON

	This work	Tang [5]	Chen [6]	Tang [7]
Max UEs	32	32	8	16
Algorithm	LAMA	CHEMP	CHEMP	EPD ^a
Soft-in soft-out	yes	no	no ^b	no
Modulation	256-QAM	256-QAM	QPSK	256-QAM
Realistic channels	yes	no	no	yes
Technology [nm]	28	40	40	28
Supply [V]	0.9	0.9	0.9	1.0
Area [mm ²]	0.37	0.58	0.076	2.0
Frequency [MHz]	400	425	500	569
Power [mW]	151	220.6	77.9	127
Throughput [Gb/s]	0.354	2.76	8	1.80
Energy ^c [pJ/b]	426	79.9	9.74	70.56
Area Eff. ^d [Gb/s/mm ²]	0.95	4.76	105.26	0.90
Norm. Energy ^{c,e,f} [pJ/b]	426	39.16	76.33	215
Norm. Area Eff. ^{d,e,f} [Gb/s/mm ²]	0.95	13.87	19.18	0.21

^aexpectation-propagation, ^bsoft-output support only, ^cenergy efficiency is power/throughput, ^darea efficiency is throughput/area, ^etechnology normalized to 28nm, $V_{dd} = 0.9$ V, ^fnormalized by $(U/32)^2$.

[2] L. Lu, G. Y. Li, A. L. Swindlehurst, A. Ashikhmin, and R. Zhang, “An Overview of Massive MIMO: Benefits and Challenges,” *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 742–758, Oct. 2014.

[3] C. Jeon, R. Ghods, A. Maleki, and C. Studer, “Optimality of large MIMO detection via approximate message passing,” in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jun. 2015, pp. 1227–1231.

[4] T. Narasimhan and A. Chockalingam, “Channel hardening-exploiting message passing (CHEMP) receiver in large-scale MIMO systems,” *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 847–860, Oct. 2014.

[5] W. Tang, C. Chen, and Z. Zhang, “A 0.58mm^2 2.76Gb/s 79.8pJ/b 256-QAM massive MIMO message-passing detector,” in *IEEE Symp. VLSI Circuits*, Jun. 2016, pp. 1–2.

[6] Y. Chen, C. Cheng, T. Tsai, W. Sun, Y. Ueng, and C. Yang, “A 501mW 7.61Gb/s integrated message-passing detector and decoder for polar-coded massive MIMO systems,” in *IEEE Symp. VLSI Circuits*, Jun. 2017, pp. C330–C331.

[7] W. Tang, H. Prabhu, L. Liu, V. Öwall, and Z. Zhang, “A 1.8Gb/s 70.6pJ/b 128×16 link-adaptive near-optimal massive MIMO detector in 28nm UTBB-FDSOI,” in *IEEE Int. Solid-State Circuits Conf. (ISSCC)*, Feb. 2018, pp. 224–226.

[8] C. Studer, S. Fateh, and D. Seethaler, “ASIC implementation of soft-in soft-out MIMO detection using MMSE parallel interference cancellation,” *IEEE J. Solid-State Circuits*, vol. 46, no. 7, pp. 1754–1765, Jul. 2011.

[9] S. Rangan, P. Schniter, and A. Fletcher, “On the convergence of approximate message passing with arbitrary matrices,” in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, June 2014, pp. 236–240.

[10] L. Cannon, “A cellular computer to implement the Kalman filter algorithm,” Ph.D. dissertation, Montana State University, USA, 1969.

[11] P. Kyösti, J. Meinilä, L. Hentilä *et al.*, “WINNER II channel models. D1.1.2 V1.2,” Tech. Rep. IST-4-027756 WINNER II, 2007.