Efficiency Metrics and Bandwidth
- A Memory Perspective

Norbert Wehn
wehn@eit.uni-kl.de
//ems.eit.uni-kl.de

Invasic Seminar
March 23 2011, Erlangen

# Content

**Computing increase and power challenge in (embedded) computing**
- Heterogeneous multi-core architectures with dedicated accelerators
- New paradigm e.g. invasive computing

**New Challenges**
- Memory and bandwidth

**Metrics for design space exploration**
- Wireless baseband processing
- Impact of memories and data transfers on metrics
- Impact of application (communications) performance on metrics

**3D MPSoCs**
- 3D memories and memory controllers

# Communication Centric World

**CAGR 2010 – 2014***

| | |
|---|---|
| Industrial | 4,9% |
| Automotive | 8,2% |
| Consumer | 4,5% |
| Wireless | 8,9% |
| Wired | 4,9% |
| Data Processing | 1,5% |

- New cellular mode is added every 3 years
- A new frequency band is added every year
- Continuous demand for higher data rates and more services

# Baseband Receiver Structure

RF&ADC → SIGNAL DETECTION → De-Interleaver → CHANNEL DECODER

SIGNAL DETECTION ↕ PARAMETER ESTIMATION

**FRONT END**   **INNER RECEIVER**   **OUTER RECEIVER**

**3.5G Digital Workload (100GOPS@1Watt)**

Future: 1 TOPS in 1+ Watt

Source: Kees van Berkel, MPSoC2010



**Mobile Phone Trends – Inner Receiver**

Source: Kees van Berkel, MPSoC2010

3

**Mobile Phone Trends – Outer receiver**

Source: Kees van Berkel, MPSoC2010

3G ⇨ 3.9G (LTE): 130 x decoding Mops/mW



**Music Baseband SDR Chip @ 65nm**

Source: Infineon

# LETI / TU KL Magali Chip

- 47■■■■■■■■■■■■■■■■■■■■■■■■■■4G SDR
- 96■■■



NOC_PERF
NoC
Evaluation

128 Mb

DCM_EXT
Ext memory
controller

External RAM
Interface
Parallel 32 bits
6.4 Gb/s

NoC Interface
Bidir 2*32 bits
2* 6.4 Gb/s

**2 processing units:**
 VLIW processors
ARM11 processor
ASIP processor
HW accelerators
**Distributed memory**
5 asynchronous
NoC router
Sophisticated power
management

Source: LETI

# Metric – Energy Efficiency

**Example - SODA, DSP and GP Architectures**

## Metric Assessment - Channel Decoders

All architectures based on standard synthesis flows, 65nm technology@worst case, all data in-house available

| Decoder | Flexibility | Max Block-size | Payload Throughput [Mbit/s] | Freq. [MHz] | Area [mm2] | Dynamic Power [mWatt] |
|---|---|---|---|---|---|---|
| ASIP (Magali) | Conv. Codes Binary TC Duo-binary TC | N=16k | 40 14(6iter) 28(6iter) | 385 (P&R) | 0.7 (P&R) | ~100 |
| LTE Turbo (Music) | LTE turbo code | N=18k | 150 (6iter) | 300 (P&R) | 2.1 (P&R) | ~300 |
| LDPC flex (Magali) | R=1/4 to R=9/10 | N=16k | 150-300 (20-10iter) | 385 (P&R) | 1.172 (P&R) | ~389 |
| LDPC fixed (Magali) | R=3/4 | N=1.2k | 480 (6iter) | 435 (P&R) | 0.583 (P&R) | ~202 |
| LDPC WiMedia 1.5 | R=1/2-4/5 | N=1.3k | 640 (R=1/2,5iter) 960 (R=3/4,5iter) | 265 | 0.51 | ~193 |
| CC Decoder | 64-state NSC | | 500 | 500 | 0.1 | ~37 |

## Algorithmic Throughput Calculations [GOPs]

| Code | Operations per decoded information bit normalized to ~8bit addition | | Infobit-Throughput ⇒Giga operations per second [GOPs] | | |
|---|---|---|---|---|---|
| | | | 100Mbit/s | 300Mbit/s | 1 Gbit/s |
| CC: states=64 | ~200 | | ~20 | ~ 60 | ~200 |
| LDPC Min-Sum (x3.4 appr. BP) | 5 iter | 75/R | ~7.5/R | ~22.5/R | ~ 75/R |
| | 10 iter | 150/R | ~15/R | ~ 45/R | ~ 150/R |
| | 20 iter | 300/R | ~ 30/R | ~ 90/R | ~ 300/R |
| | 40 iter | 600/R | ~ 60/R | ~ 180/R | ~ 600/R |
| Turbo Max-Log | 2 iter | 280 | ~ 28 | ~ 84 | ~ 280 |
| | 4 iter | 560 | ~ 56 | ~168 | ~ 560 |
| | 6 iter | 840 | ~ 84 | ~252 | ~ 840 |

6

## Area- and Energy Efficiency



---

## What about Memory/Data Transfers

Current metric: energy efficiency = only operations/energy

Data transfers/ accesses substantially contribute to the power consumption

Example (R=0.5)

150 Mbit/s Turbo : ~126 Gops     ~40 Gaccesses

150 Mbit/s LDPC : ~90 Gops     ~80 Gaccesses

Efficient data transfer is key for efficient implementation

- LTE TC: special interleaver structure to avoid access conflicts
- DVB-S2/WiMAX LDPC: special code structure to minimize access conflicts

Efficiency metrics based on operations only are not appropriate

- Power includes operations and accesses!
- Architectures are favored where operations dominate compared to accesses

# Decoders in System Design Space



Chart — Energy Efficiency: decoded bit/energy (bit/nJ) vs Area Efficiency: (Mbit/s)/mm²

Legend:
- ASIP TC (14Mbit/s)
- LTE TC (150Mbit/s)
- LDPC flexible (~100 Mbit/s)
- LDPC WiMedia (~1Gbit/s)
- CC (500Mbit/s)

Labels: Min-Sum, Λ-3-Min

# Communications Performance

**Overall efficiency of a baseband receiver depends on**

- **Implementation performance**
- **Communications performance**
- **Flexibility**

**Scenario 1: Fixed Communication performance**

- **Comparison of two iterative decoders with same communications performance but different parameters (codes, code rate, iterations)**
- ⇨ impact on implementation efficiency

**Scenario 2: Implementation driven**

- **Comparison of iterative and non-iterative decoders with varying communications performance**
- **64-state convolutional code 960 Mbit/s (WiMedia 1.2) and WiMedia 1.5 LDPC decoder**
- ⇨ impact on implementations efficiency

## Scenario 1: Fixed Communication Performance



Blocksize 6145 Information bits, TC: 150Mbit/s @6.5 iterations

## Scenario 1: Implementation Efficiency

## Scenario 2: Varying Communication Performance



FER vs $E_B/N_0$ [dB]

Legend:
- LDPC, R=3/4, 5 iterations
- LDPC, R=3/4, 2 iterations
- CC, R=3/4
- LDPC, R=3/4, 1 iterations

**Simulation Set-Up: WiMedia 1.5 chain, 16-QAM, Blocksize 1200 Information bits**

## Scenario 2: Implementation Efficiency



Energy Efficiency: decoded bit/energy (bit/nJ)

- 4dB worse com. Perf. (1iter)
  Identical Throughput
- 4dB worse com. Perf. (1iter)
  5 times higher throughput
- Identical comm. Perf. (2iter)
  Identical throughput
- Identical com. Perf. (2iter)
  2.5 times higher throughput
- 4dB better com. Perf. (5iter)
  Identical throughput

▲ LDPC WiMedia1.5
● CC WiMedia1.5

Area Efficiency: (Mbit/s)/mm2

10

# Lessons learned

- Understanding trade-offs between implementation efficiency, application performance and flexibility requirements is mandatory for efficient baseband receivers

- Operation based metrics for energy and area efficiency can be misleading

- Memory and data transfers have to be considered in metrics for design space exploration

- Implementation efficiency metrics have to be linked to application performance ⇨ trajectory

# Off-chip Memory Bandwidth



**Using Cache size to accommodate increasing traffic is VERY expensive!**

$C \approx (T/B)^{\delta}$
with $\delta$..2-3

2x increased traffic drives 8x cache size
(constant memory bandwidth)

4x increased traffic drives 64x cache size
(constant memory bandwidth B)

Source: IBM

## Next-Generation Mobile Platform Traffic

| Graphic | Video | Display | Base-band |

NoC

Multi-ported front-end

Scheduler

>>10GBps

DRAM Controller

| Channel 1 DDR3 | Channel 2 DDR3 |

Memory Architecture

■ **Traditional JEDEC DRAM channels are saturating**

## Next Generation Teraflop Computing Platform

Year 2018
8nm Core, 10Gflop

400mm² Die size
1150 Cores

DP FP Add, Multiply
Integer Core, RF
Router
0.17mm² (50%)

Memory 0.35MB
0.17mm² (50%)

~0.6mm

20mm

400mm2

128 GB    128 GB

256GB/s    64b

128 GB    128 GB

Source: Intel

1 TF, ~ 100 W

Aggressive
voltage
scaling

Compute

Fabric

Hierarchical
heterogeneous
topologies

DRAM

Efficient signaling
Repartitioning

**Goal of 20W**

# 3D Integration with TSVs



**Cu Pillars**

**Top Die**

**TSV**

**Silicon Interposer**

**Substrate**

**Through Silicon Vias (TSV)**
- **Polysilicon filled (FEOL)**
  - 10.000 TSV/mm$^2$
- **Copper filled (BEOL)**
  - 500 TSV/mm$^2$

High-Bandwidth, Low-Latency Connections

Microbumps

Through Silicon Vias (TSV)

C4 Bumps

28 nm FPGA Die Slices

Silicon Interposer

Package Substrate

BGA Solder Balls

Source: XILINX

3D TSV

DRAM

NVM

Discrete passives

RRAM

BEOL RRAM

SoC

Source: LETI

# Wide IO Technology (JEDEC Standard 2012)

**Channel**
- 4 x 64Mb
- 128 bit @ 200MHz SDR
- 3.2GBps

**Memory**
- 4 channels
- 1Gb
- 512 bit IO
- 12.8GBps



Channel 0   Channel 1

Bank 0 | Bank 1   Bank 0 | Bank 1
Bank 2 | Bank n   Bank 2 | Bank n

Interconnect area

Bank 0 | Bank 1   Bank 0 | Bank 1
Bank 2 | Bank n   Bank 2 | Bank n

Channel 2   Channel 3

Wide-IO DRAM

Mobile   Processor

Array of ~1000 microbumps, 40/50μm pitch

Array of TSVs 10μm Ø, 40μm pitch

Source: LETI

## Power Savings in DRAM Memory Interfaces

■ **Much wider I/Os possible >> 32 bits**

| Memory link, peak bandwidth and power consumption efficiency | | | Cost for 1TBps memory bandwidth | |
|---|---|---|---|---|
| | | | Number of data IO pins | Interface power consumption |
| **Computing memory IF standard** | Multi-core SoC → DDR3 → DRAM | 8.532 GBps 30 mW/Gbps | 3800 | 240 W |
| | 1066 MHz I/O bus clock, 32 bits, 1.5 V, Double Data Rate | | | |
| **Mobile memory IF standards** | Multi-core SoC → LPDDR2 → DRAM | 4.264 GBps 20 mW/Gbps | 7700 | 160 W |
| | 533 MHz I/O bus clock, 32 bits, 1.2 V, Double Data Rate | | | |
| | Multi-core SoC → Wide I/O → DRAM | 12.8 GBps 4 mW/Gbps | 41000 | 32 W |
| | 200 MHz I/O bus clock, 512 bits, 1.2 V, Single Data Rate | | | |

## Transition to 3D-DRAMs



14

## Organization and Naming

**A single 3D-layer consists of 3D-DRAM banks**

**A bank is composed of DRAM core tiles**

3D-layer (= tier)

3D-DRAM bank

64Mb   64Mb

3D-DRAM core tile

64M ARRAY
ROW
COLUMN
Control / Power generators / Signaling
TSVs Power & Signals

Tile: basic memory macro cut out of a commodity DRAM

---

## Investigated Technologies

**Example: 64Mb 3D-DRAM core tile**

- TSV areas added
- Deep trench / buried WL / Stack
- Cell sizes: $8F^2 - 4F^2$
- Based on measured* & simulated data

64M ARRAY
ROW
COLUMN
Control / Power generators / Signaling
TSVs Power & Signals

| No. | Techn. node | Cell size | Cell type | Area [mm²] | Row $t_{RAS}$ [ns] | Row -> Col. $t_{RCD}$ [ns] | Column $t_{CCD}$ [ns] |
|---|---|---|---|---|---|---|---|
| 1 | 75nm | $8F^2$ | Deep Trench | 5.20 | 39.0 | 9.30 | 6.05 |
| 2 | 65nm | $6F^2$ | buried WL | 3.54 | 27.1 | 7.45 | 5.42 |
| 3 | 58nm | $6F^2$ | Stack | 3.00 | 31.9 | 7.31 | 4.70 |
| 4 | 46nm | $6F^2$ | buried WL | 2.26 | 26.4 | 6.44 | 3.59 |
| 5 | 45nm | $4F^2$ | buried WL | 1.92 | 26.0 | 5.98 | 2.76 |

# Single 3D-layer Design Space

2x 64Mb with 128 I/Os

1x 128Mb with 64 I/Os

**128Mb**  ROW

**64Mb**  **64Mb**

**COLUMN**

Control / Voltage generators / Signaling

TSV area: I/Os **64** and Power

4x 32Mb with 256 I/Os

**32Mb**  **32Mb**

**32Mb**  **32Mb**

8x 16Mb with 512 I/Os

16Mb  16Mb 16Mb  16Mb

16Mb  16Mb 16Mb  16Mb

---

# 3D-DRAM models

**Inputs:**
- Number of 3D **layers**
- Number of log. **banks**
- DRAM **size** (Mbit)
- Data **IO width**
- **IO width** per base
- DDR or SDR interface
- Long or short bitlines
- Long or short wordlines
- TSV pitch and diameter
- Number of power TSVs
- **Technology node**

**3D-DRAM generator model**

Layer 3
Layer 2
Layer 1
Layer 0

**Outputs:**

**Timings:**
- $t_{RCD}$, $t_{RAS}$, $t_{CCD}$,
- $t_{RC}$, $t_{RP}$, $t_{WR}$, …
- Max. frequency

**Power values:**
- Active (RD/WR)
- Standby
- Power down
- Self refresh

**Area:**
- Per 3D layer

Cycle-accurate SystemVerilog /
SystemC simulation models

## Metrics for Exploration

- **Throughput (TP)**
  - maximal theoretical bandwidth ($f_{max} \cdot$ IO width)
  - $f_{max}$ determined by architecture & technology
    <u>here</u>: column to column access delay ($t_{CCD}$)

- **Area efficiency**
  - Maximum learning out of the commodity DRAM production: minimize cost/bit
  - Maximize cell efficiency (CE) = memory cell area / total area [%]

- **Energy efficiency (EE)**
  - TP / average power = access / energy [MB/mJ]

---

## Single 3D-layer design space

- Cell efficiency vs. max. theoretical throughput (TP) for various banks



(a)  1x 128Mb
(b)  2x 64Mb
(c)  4x 32Mb
(d)  8x 16Mb

# Single 3D-layer design space

- Energy efficiency (EE) vs. Cell efficiency (CE)



Single 3D-layer ⇨ bank composed of 2x 64Mb tiles
independent of technology

# Comparison 1Gb

- 1Gb, 8 bank standard 2D-DRAM wo/ IO driver and termination power
- 1Gb stacked eDRAM: extrapolated published 2,39Mb SOI macro [ISSCC]
- 1Gb Mobile Low-Power DDR SDRAM x 16 wo/ IO driver and termination



8x 128Mb = 1Gbit

## Multi-Channel 3D-DRAM Controller

**Different request granularities**

- **But normally fixed to 128 bit (Wide IO JEDEC Standard)**

ARM-MP Core | Video Core | GPU Core

32-bit | 64-bit | 128-bit

**3D-DRAM Memory controller**

**Three types of I/O accesses possible** — **from 32-bit to 128-bit**

**3D-DRAM**

---

## Multi-Channel 3D-DRAM Controller

**Front End:**

- Synchronization with Dual Clock FIFOs
- Arbitration
- Buffering, Scheduling, Reordering

**Back End**

- 3D DRAM command Encoding
- Tracking of the BANK status
- Multi IO reconfiguration and data latching for 32/64/128 bit

FE Memory Controller      Channel Controller

32

64

128

Memory Controller Front End

cc3
cc2
cc1
cc0

# Fine grained 3D-DRAM access

On the fly switching: 32, 64, 128

**32 Bit**

**64 Bit**

**128 Bit**

Memory Controller

Stacked 3D DRAM

**1** Wordline **8** CSLs

**1** Wordline **8** CSLs

**1** Wordline **4** CSLs

**1** Wordline **8** CSLs

128Mb  ROW  ROW  128Mb

COL  COL

ctrl  ctrl

IOs  IOs

64  64

32  64  128

**128bit**   **Native – 16 CSLs – 2 Wordlines**

**64bit**   **8 CSLs + 1 Wordline**

**32bit**   **4 CSLs + 1 Wordline**

---

# Flexible 3D-DRAM system

DRAM 0
DRAM 1
DRAM 2
DRAM 3
DRAM 4
DRAM 5
DRAM 6
DRAM 7

3D-DRAM cube

CC Controller Frontend

CCD

Logic Chip

**Vertical Channel**   Proximity Channel Controllers (CC)

128Mb  128Mb  64Mb  64Mb

64Mb  64Mb

Bank 7

Bank 0

**Single Channel with 128 IO's**
**Each layer: 4 x 64Mb or 2 x 128Mb**
**8 layers ⇒ 2Gb/channel**

# Investigated 3D-DRAM Configurations

| 3D-DRAM SINGLE CHANNEL CONFIGURATIONS | | | | | | |
|---|---|---|---|---|---|---|
| Dens. [Mb] | Architecture # lay. x [org.] | # of banks | Techn. [nm] | Cell size | $A_{total}$ [mm$^2$] | Freq. [MHz] |
| **SDR x128** | | | | | | |
| **256 | 1 x [4x64Mb] | 4 | 58 | 6F$^2$ | 16 | 200 |
| 512 | 2 x [4x64Mb] | 4 | 58 | 6F$^2$ | 26 | 200 |
| 1024 | 8 x [2x64Mb] | 8 | 46 | 6F$^2$ | 35 | 300 |
| *2048 | 8 x [2x128Mb] | 8 | 46 | 6F$^2$ | 60 | 167 |
| 4096 | 8 x [4x128Mb] | 8 | 45 | 4F$^2$ | 97 | 200 |
| **DDR x128** | | | | | | |
| 256 | 1 x [4x64Mb] | 4 | 58 | 6F$^2$ | 22 | 200 |
| 512 | 2 x [4x64Mb] | 4 | 58 | 6F$^2$ | 32 | 200 |
| 1024 | 8 x [2x64Mb] | 8 | 46 | 6F$^2$ | 44 | 300 |
| *2048 | 8 x [4x64Mb] | 8 | 46 | 6F$^2$ | 69 | 300 |
| 4096 | 8 x [4x128Mb] | 8 | 45 | 4F$^2$ | 98 | 200 |

** *Density emulates the published Samsung 1Gb WIDE IO chip [15].*

# Simulation Set-Up

**Emulation of workload via 3 traffic generators**

- Traffic A – Cache misses of a 2 core ARM ~ 100 MB/s per core
- Traffic B – DMA accesses in a SoC (Imaging) ~ 0.8 GB/s
- Traffic C – HD Video DMA accesses ~ 1.5 GB/s

## Results with Page hit rate of 50%



## The Future - 3D Magali Chip

- 65nm tech, 72mm$^2$, 1980TSVs for 3D NoC, 1250 TSV for wide I/O memory
- Heater, temperature sensors



Source: LETI

# Conclusion

- **Bandwidth and memory will be big challenges in future computing systems**

- **We will see new memory devices e.g. memristor based (RRAMs) or spin based memories (MRAMs)**

- **The future in computation will be 3D**

- **New heterogeneous memory architectures**

- **Large opportunity for research**