# CICC

IEEE Custom Integrated Circuits Conference

## 9-2: An Energy-Efficient and Runtime-Reconfigurable FPGA-Based Accelerator for Robotic Localization Systems

Qiang Liu[1,]*, **Zishen Wan**[2,]*, Bo Yu[3,]*, Weizhuang Liu[1], Shaoshan Liu[3], Arijit Raychowdhury[2]

* Equally Contributed Authors

[1] Tianjin University, China      [2] **Georgia Institute of Technology, USA**      [3] PerceptIn, USA

April 25, 2022

SSCS IEEE SOLID-STATE CIRCUITS SOCIETY™ | ◆IEEE

# Bio



**Email**: zishenwan@gatech.edu

**Homepage**: https://zishenwan.github.io

- **Speaker: Zishen Wan**
  - **PhD Student in Georgia Tech (20Fall-Now)**
    - **Advisor: Prof. Arijit Raychowdhury**
  - MS in Harvard University
    - Advisor: Prof. Vijay Janapa Reddi
  - BS in Harbin Institute of Technology

- **Research Interest**
  - VLSI, computer architecture, edge computing.
  - Efficient and resilient hardware and system design for autonomous machines.

# Motivation: Autonomous Systems

Drones

Self-Driving Cars

Robots

# Motivation: Autonomous Systems

Drones

Self-Driving Cars
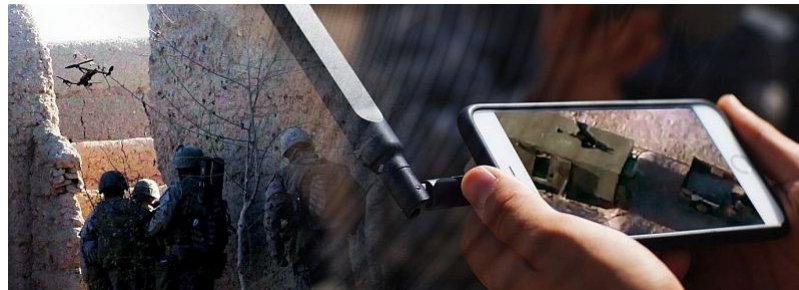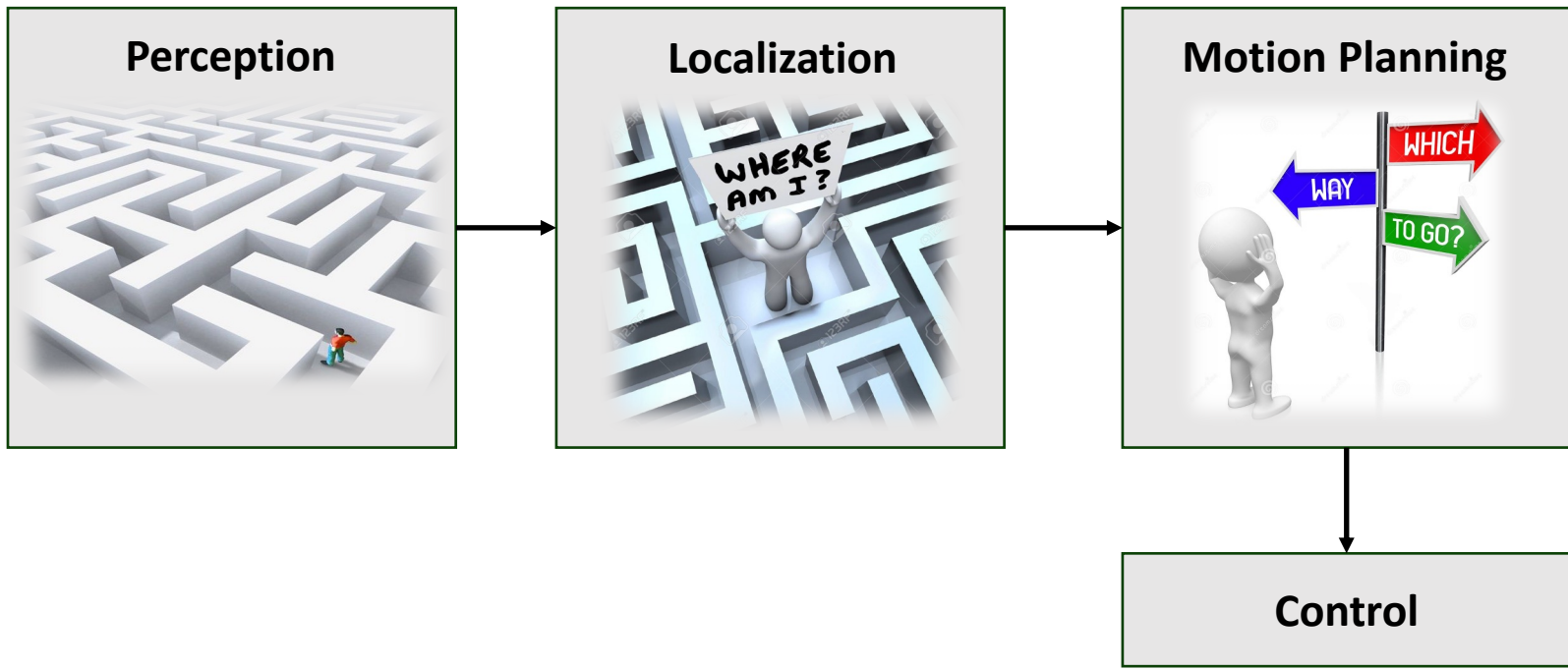
Robots

**Applications**

Search & Rescue

Package Delivery

Surveillance

# How Does Autonomous System Work?

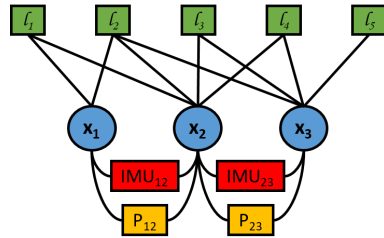| Perception | Localization | Motion Planning |
|:---:|:---:|:---:|

**Control**

**CICC**

# How Does Autonomous System Work?

# Challenges

**Large Factor Graph:**



*4000+*

*factors*

# Challenges

## Large Factor Graph:



*4000+ factors*

## Real-Time Requirement:

# Challenges

**Large Factor Graph:**



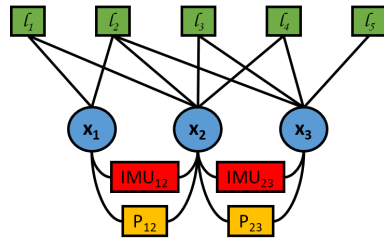***4000+ factors***

**Real-Time Requirement:**



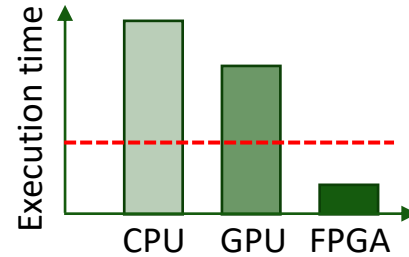**Low Power Budget:**



Big battery
CPU/GPU: 10-100W

# Challenges

## Large Factor Graph:



*4000+ factors*

## Real-Time Requirement:



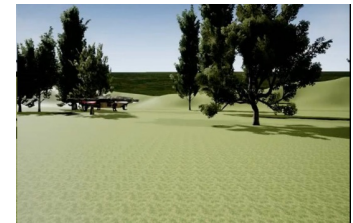## Low Power Budget:



Big battery
CPU/GPU: 10-100W

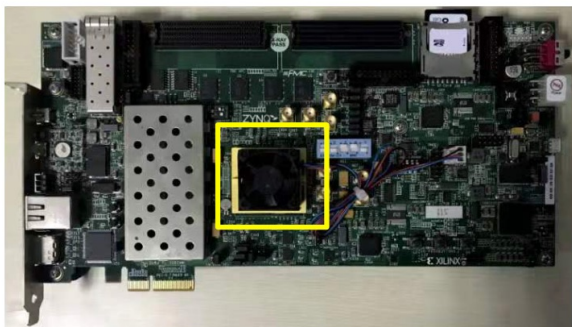## Dynamic Changing Environments:



Sparse                    Dense

# Energy-Efficient Localization and Mapping



FPGA Zynq-7000 SoC ZC706
with XC7Z045 FFG900-2

- Energy-efficient & real-time localization and mapping

- Dynamic reconfiguration at runtime

- Real-time performance of 61 fps at 3.45W (56mJ/frame)
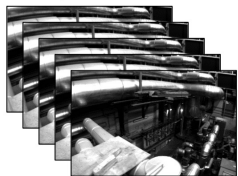
# Outline

- SLAM: Simultaneously Localization & Mapping

- Hardware Architecture

- Main Contributions

- Evaluations and Comparisons

- Summary

# Outline

- SLAM: Simultaneously Localization & Mapping
- Hardware Architecture
- Main Contributions
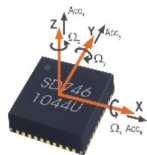- Evaluations and Comparisons
- Summary
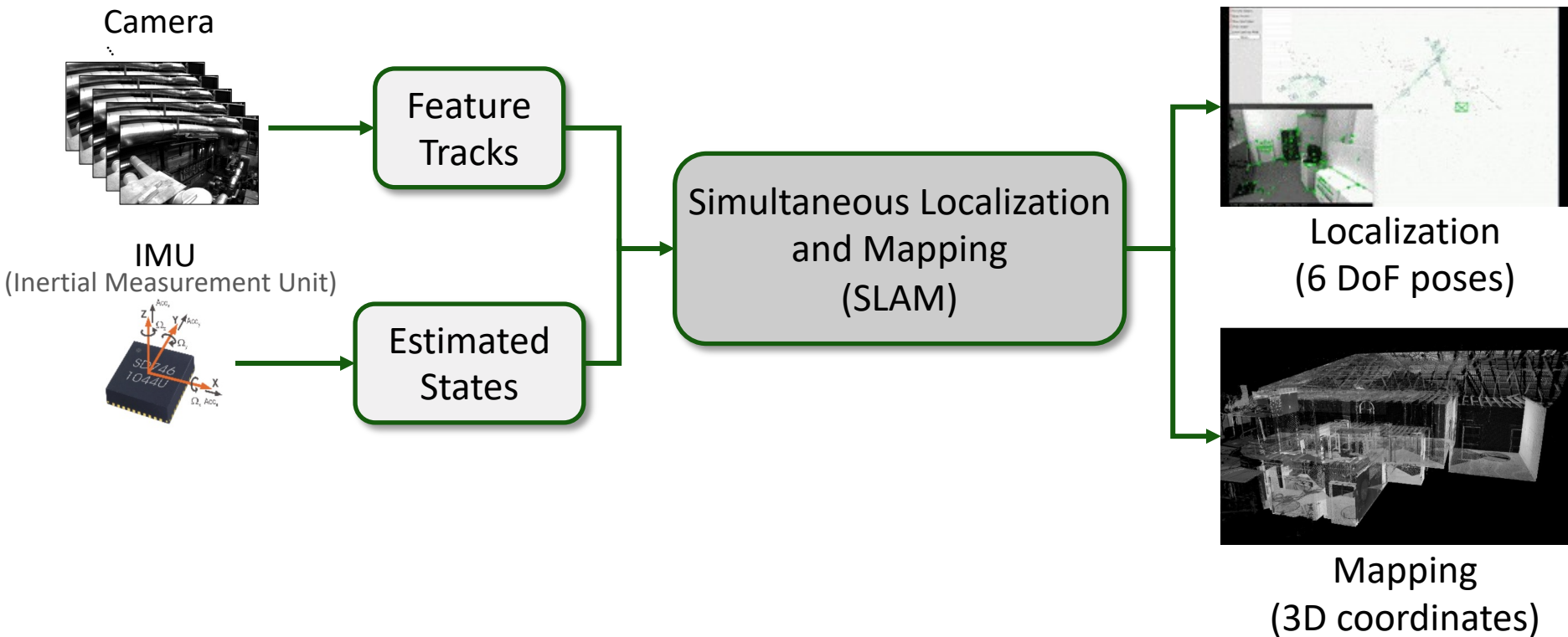
# Localization and Mapping Using SLAM

Camera



Feature Tracks

IMU
(Inertial Measurement Unit)



Estimated States

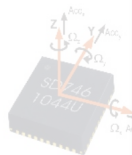# Localization and Mapping Using SLAM

Camera

Feature Tracks

IMU
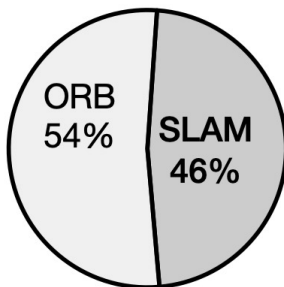(Inertial Measurement Unit)

Estimated States

Simultaneous Localization and Mapping (SLAM)

Localization
(6 DoF poses)

Mapping
(3D coordinates)

# Localization and Mapping Using SLAM

Camera

IMU
(Inertial Measure...

**SLAM is computationally intensive:**
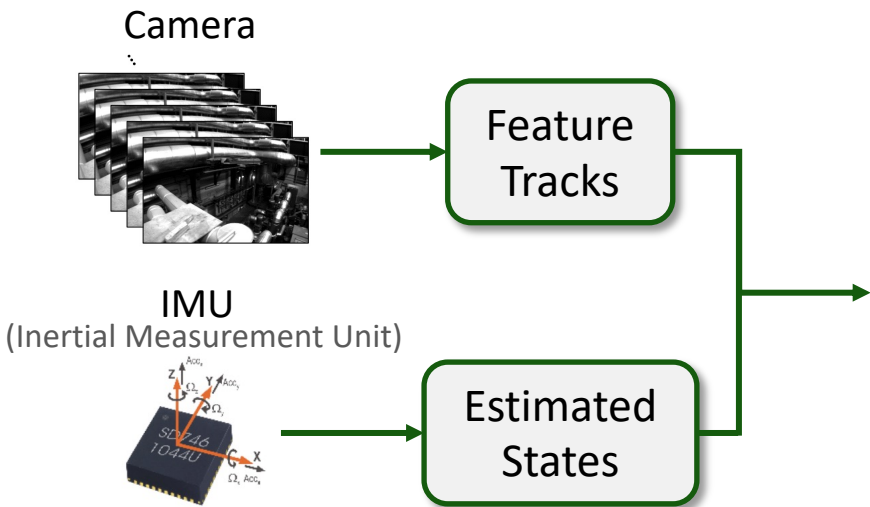
**ORB-SLAM**

FrontEnd: ORB
BackEnd: SLAM

ORB 54%  SLAM 46%

**LK-SLAM**

FrontEnd: LK
BackEnd: SLAM

LK 22%  SLAM 78%

...ization
(... poses)

Mapping
(3D coordinates)

# How Does SLAM Work?

Camera

Feature Tracks

IMU
(Inertial Measurement Unit)

Estimated States

# How Does SLAM Work?

Camera

Feature Tracks

IMU
(Inertial Measurement Unit)

Estimated States

**ORB-SLAM**

FrontEnd: ORB
BackEnd: SLAM

ORB 54% | SLAM 46%

**LK-SLAM**

FrontEnd: LK
BackEnd: SLAM

LK 22% | SLAM 78%

# How Does SLAM Work?

Camera



**Feature Tracks**

IMU
(Inertial Measurement Unit)



**Estimated States**

**Maximum a Posteriori (MAP) Estimation**

Nonlinear least squares (NLS) optimization problem:

$$\min_{\mathbf{p}}\{\sum_{i=1}^{N}||\mathbf{r}_p - \mathbf{H}_p\mathbf{p}||^2 + ||\mathbf{o}_i - \mathcal{P}_i(\mathbf{p})||^2_{C_i}\}$$

State

**ORB-SLAM**

FrontEnd: ORB
BackEnd: SLAM

ORB 54% | SLAM 46%

**LK-SLAM**

FrontEnd: LK
BackEnd: SLAM

LK 22% | SLAM 78%

**State:**
6 DoF poses      (location)
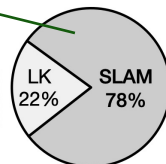3D coordinates (mapping)

# How Does SLAM Work?

Camera

IMU
(Inertial Measurement Unit)

Feature Tracks

Estimated States

Marginalization

Maximum a Posteriori (MAP) Estimation

Nonlinear least squares (NLS) optimization problem:
$$\min_{\mathbf{p}}\{\sum_{i=1}^{N} ||\mathbf{r}_p - \mathbf{H}_p\mathbf{p}||^2 + ||\mathbf{o}_i - \mathcal{P}_i(\mathbf{p})||_{C_i}^2\}$$

State

**State:**
6 DoF poses (location)
3D coordinates (mapping)

**ORB-SLAM**

FrontEnd: ORB
BackEnd: SLAM

ORB 54%  SLAM 46%

**LK-SLAM**

FrontEnd: LK
BackEnd: SLAM

LK 22%  SLAM 78%

# How Does SLAM Work?



Camera

IMU
(Inertial Measurement Unit)

Feature Tracks

Estimated States

Marginalization

Maximum a Posteriori (MAP) Estimation

State

Nonlinear least squares (NLS) optimization problem:

$$\min_{\mathbf{p}}\{\sum_{i=1}^{N} ||\mathbf{r}_p - \mathbf{H}_p\mathbf{p}||^2 + ||\mathbf{o}_i - \mathcal{P}_i(\mathbf{p})||_{C_i}^2\}$$

**ORB-SLAM**

FrontEnd: ORB
BackEnd: SLAM

ORB 54% | SLAM 46%

**LK-SLAM**

FrontEnd: LK
BackEnd: SLAM

LK 22% | SLAM 78%

Jacobian Matrix

Schur Elimination

Cholesky Decomposition

**State:**
6 DoF poses      (location)
3D coordinates (mapping)

# How Does SLAM Work?



SLAM Latency Distribution

Cholesky Decomposition

NLS Solver 50.15%

Schur Elimination

Jocobian & Residual

Marginalization 44.40%

Others 10.80%

ORB-SLAM
FrontEnd: ORB
BackEnd: SLAM

ORB 54% | SLAM 46%

LK-SLAM
FrontEnd: LK
BackEnd: SLAM

LK 22% | SLAM 78%

Marginalization

Maximum a Posteriori (MAP) Estimation

Nonlinear least squares (NLS) optimization problem:

$$\min_{\mathbf{p}}\{\sum_{i=1}^{N}||\mathbf{r}_p - \mathbf{H}_p\mathbf{p}||^2 + ||\mathbf{o}_i - \mathcal{P}_i(\mathbf{p})||^2_{C_i}\}$$

State

Jacobian Matrix | Schur Elimination | Cholesky Decomposition

**State:**
6 DoF poses    (location)
3D coordinates (mapping)

# Outline

- SLAM: Simultaneously Localization & Mapping
- Hardware Architecture
- Main Contributions
- Evaluations and Comparisons
- Summary

# Hardware Architecture - Overview

# Hardware Architecture – Perception



**Sensor Input:**
Camera + IMU, process in host

# Hardware Architecture – SLAM (NLS Optimization)



**SLAM Nonlinear Least Squares (NLS) Optimization:**
Jacobian, Schur elimination, Cholesky Decomposition, etc

# Hardware Architecture – SLAM Marginalization



**SLAM Marginalization:**
Jacobian, Schur elimination, Cholesky Decomposition, etc

# Outline

- SLAM: Simultaneously Localization & Mapping

- Hardware Architecture

- Main Contributions

- Evaluations and Comparisons

- Summary

# Method 1

Data Reuse

# Data Reuse & Design Hierarchy



2 **K**eyframes
3 **F**eature Points (F1~F3)
4 **O**bservations (O1~O4)

# Data Reuse & Design Hierarchy



Jacobian Matrix

2 **K**eyframes
3 **F**eature Points (F1~F3)
4 **O**bservations (O1~O4)

<feature point, observation>
pairs have non-zero values

# Data Reuse & Design Hierarchy



F1 F2 F3

O1
O2
... O3 O4

Keyframe 1    Keyframe 2

2 **K**eyframes
3 **F**eature Points (F1~F3)
4 **O**bservations (O1~O4)

O1  O2  O3  O4
F1
F2
F3

Jacobian Matrix

<feature point, observation>
pairs have non-zero values

Keyframes
RAM → Keyframe Block → Rotation matrix RAM

FIFO
Observations → Observation Block → Jacobian matrix

FIFO
Feature points → Feature Block → FIFO Point coordinates

# Data Reuse & Design Hierarchy



2 **K**eyframes
3 **F**eature Points (F1~F3)
4 **O**bservations (O1~O4)

Jacobian Matrix

<feature point, observation>
pairs have non-zero values

**Three-Level Block Designs**:
- Keyframe-level:     Rotation matrix of keyframes
- Feature-level:      3D coordinates
- Observation-level: Jacobian matrix

# Data Reuse & Design Hierarchy



F1  F2  F3

O1  O2  O3  O4

O1
O2
O3  O4

Keyframe 1    Keyframe 2

2 **K**eyframes
3 **F**eature Points (F1~F3)
4 **O**bservations (O1~O4)

Jacobian Matrix

<feature point, observation>
pairs have non-zero values

Keyframes
RAM → **Keyframe Block** → Rotation matrix RAM

FIFO Observations → **Observation Block** → Jacobian matrix

FIFO Feature points → **Feature Block** → FIFO Point coordinates

**Two-Level Data Reuses**:
- Feature-reuse: across associated observations
- Keyframe-reuse: over all obsn. within keyframe

**Three-Level Block Designs**:
- Keyframe-level:    Rotation matrix of keyframes
- Feature-level:     3D coordinates
- Observation-level: Jacobian matrix

# Data Reuse & Design Hierarchy



2 **K**eyframes
3 **F**eature Points (F1~F3)
4 **O**bservations (O1~O4)

Jacobian Matrix

<feature point, observation>
pairs have non-zero values

**Two-Level Data Reuses**:
- Feature-reuse: across associated observations
  ➡ feature (row)-stationary
- Keyframe-reuse: over all obsn. within keyframe

**Three-Level Block Designs**:
- Keyframe-level:     Rotation matrix of keyframes
- Feature-level:      3D coordinates
- Observation-level: Jacobian matrix

# Method 2

## Symmetry & Sparsity

# Diagonal Computation + Symmetry + Hardware Reuse

Shure Elimination:

**A**     x    $\Delta$**p**   =   **b**

8

8

x

=

# Diagonal Computation + Symmetry + Hardware Reuse

Shure Elimination:



A →
$$\begin{bmatrix} U & X \\ W & V \end{bmatrix}$$

# Diagonal Computation + Symmetry + Hardware Reuse

Shure Elimination:



**Make U as diagonal matrix:**
$O(n^3)->O(n)$ computational complexity

**X becomes the transpose of W:**
    **1.34x** on-chip memory reduction

# Diagonal Computation + Symmetry + Hardware Reuse

Shure Elimination:



$$A \rightarrow \begin{bmatrix} U & X \\ W & V \end{bmatrix}$$

Marginalization:



Goal: prior information $H_p$

$$H_p = A - ZM^{-1}Z^T$$

$M$ is general matrix

$$M = \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix}$$

$$M^{-1} = f(M, S')$$

$$S' = M_{22} - M_{21}M_{11}^{-1}M_{12}$$

Make $M_{11}$ diagonal

**Reuse** Schur Elimination and Cholesky Decomposition

# Diagonal Computation + Symmetry + Hardware Reuse

Shure Elimination:



Marginalization:

# Diagonal Computation + Symmetry + Hardware Reuse

**<u>Make M as diagonal matrix:</u>**
$O(n^3) \rightarrow O(n)$ computational complexity

**<u>Reuse Schur Elimination circuit in Marginalization:</u>**
Reduce resource consumption without performance degradation

Marginalization:



Information matrix H

6    4

6    M [6x6]    $Z^T$ [6x4]
     $M_{11}$    $M_{12}$
     $M_{21}$    $M_{22}$

4    Z [4x6]    A [4x4]

Goal: prior information $H_p$
$$H_p = A - ZM^{-1}Z^T$$
M is general matrix

$$M = \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix}$$

$$M^{-1} = f(M, S')$$

$$S' = M_{22} - M_{21}M_{11}^{-1}M_{12}$$

Make $M_{11}$ diagonal

**Reuse** Schur Elimination and Cholesky Decomposition

# Data Layout + Symmetry + Sparsity

**S matrix**



**S** matrix: store the parameters
for the system
(40%-80% of total storage)

720 kb

# Data Layout + Symmetry + Sparsity

**S matrix**

S matrix: store the parameters
for the linear system
(40%-80% of total storage)

720 kb

15     IMU

+

6     Vision

# Data Layout + Symmetry + Sparsity

**S matrix**



**S** matrix: store the parameters
for the linear system
(40%-80% of total storage)

720 kb

# Data Layout + Symmetry + Sparsity

**S matrix**

**S** matrix: store the parameters
for the linear system
(40%-80% of total storage)

720 kb

15    IMU
Symmetry

+

6    Vision
Symmetry

# Data Layout + Symmetry + Sparsity



**S matrix**

**S** matrix: store the parameters
for the linear system
(40%-80% of total storage)

15    IMU

Symmetry

+

6    Vision

Symmetry    Co-observation

4.1x reduction

720 kb ———————————————→ 175.97 kb

# Data Layout + Symmetry + Sparsity



**Data Layout + Symmetry + Sparsity + Co-observation**

**4.1x** memory reduction

Exploiting data characteristics unique to SLAM

S matrix: store the parameters
for the linear system
(40%-80% of total storage)

720 kb ⟶ 175.97 kb

4.1x reduction

# Method 3

## Time-Multiplex & Pipeline

# Time-Multiplexed + Pipeline Processing

Cholesky decomposition: $S = LL^T$ ($S$: symmetric matrix; $L$: lower triangular matrix)

# Time-Multiplexed + Pipeline Processing

Cholesky decomposition: $S = LL^T$ ($S$: symmetric matrix; $L$: lower triangular matrix)

# Time-Multiplexed + Pipeline Processing

Cholesky decomposition: $S = LL^T$ ($S$: symmetric matrix; $L$: lower triangular matrix)

# Time-Multiplexed + Pipeline Processing

Cholesky decomposition: $S = LL^T$ ($S$: symmetric matrix; $L$: lower triangular matrix)



**Time-Multiplexed + Pipeline:**

**3.3x** Hardware resources reduction
**5.75x** Processing time reduction

# Method 4

Runtime Reconfiguration & Clock Gating

# Runtime Reconfiguration + Clock Gating

# Runtime Reconfiguration + Clock Gating



Software Processing: Feature Points → Levenberg-Marquardt (LM) Algorithm → Marginalization Calculation → 6 DoF poses + 3D coordinantes

Hardware Operation: Sensors → NLS Solver Accelerator → Marginalization Accelerator → State Vector (Localization + Mapping)

Time

#Feature points↓ ➡ Accuracy↓ ➡ NEED #iterations↑

KITTI Dataset: 5.47W — Baseline

EuRoC Dataset: 5.48W — Baseline

# Runtime Reconfiguration + Clock Gating



**Time**

**Software Processing**

Feature Points → Levenberg-Marquardt (LM) Algorithm → Marginalization Calculation → 6 DoF poses + 3D coordinantes

**Hardware Operation**

Sensors → NLS Solver Accelerator → Marginalization Accelerator → State Vector (Localization + Mapping)

**Asynchronous**

Feature Points →

| Feature Points | # Iterations in NLS | # Schur blocks | # Update blocks |
|---|---|---|---|
| 0-200 | 6 | 47 | 97 |
| 200-250 | 5 | 42 | 63 |
| 250-300 | 4 | 35 | 42 |
| ... | ... | ... | ... |

Lookup Table (0.3 kb)

↻ Automated Self-Update with New Environments

**KITTI Dataset**

Power

5.47W

Baseline

**EuRoC Dataset**

Power

5.48W

Baseline

# Runtime Reconfiguration + Clock Gating



**Time**

**Software Processing:** Feature Points → Levenberg-Marquardt (LM) Algorithm → Marginalization Calculation → 6 DoF poses + 3D coordinantes

**Hardware Operation:** Sensors | *Runtime Reconfig. + Clock Gating (RR + CG)* | NLS Solver Accelerator | *Runtime Reconfig. + Clock Gating (RR + CG)* | Marginalization Accelerator | State Vector (Localization + Mapping)

**Asynchronous:** Feature Points →

| Feature Points | # Iterations in NLS | # Schur blocks | # Update blocks |
|---|---|---|---|
| 0-200 | 6 | 47 | 97 |
| 200-250 | 5 | 42 | 63 |
| 250-300 | 4 | 35 | 42 |
| ... | ... | ... | ... |

Lookup Table (0.3 kb)

Automated Self-Update with New Environments

**KITTI Dataset**
Power
5.47W   3.73W
1.47x
Baseline   RR+CG

**EuRoC Dataset**
Power
5.48W   3.45W
1.59x
Baseline   RR+CG

# Runtime Reconfiguration + Clock Gating

**Runtime Reconfigurable + Clock Gating:**

**1.47x** power reduction in KITTI dataset
**5.75x** power reduction in EuRoC dataset
**<0.01cm** accuracy degradation

# Outline

- SLAM: Simultaneously Localization & Mapping
- Hardware Architecture
- Main Contributions
- Evaluations and Comparisons
- Summary

# Evaluation - Dataset

- EuRoC Dataset (for drone)
  - A very challenging, and widely used UAV dataset
  - 11 sequences with three categories: easy, medium & difficult
  - This work: Machine Hall sequences

  

- KITTI Dataset (for self-driving car)
  - A widely used autonomous driving vision benchmark
  - Task of interest: stereo, optical flow, visual odometry, 3D object detection and 3D tracking
  - This work: odometry (grayscale sequence)

# Evaluation – FPGA Platform



FPGA Zynq-7000 SoC ZC706
with XC7Z045 FFG900-2

| Operation Frequency | 143 MHz |
|---|---|
| LUT | 144108 (65.92%) |
| Flip-Flop | 172935 (39.56%) |
| BRAM | 268 (49.17%) |
| DSP | 869 (96.56%) |

# Evaluation

*- Processing Latency and Energy of FPGA, CPU, and GPU*



- FPGA: Xilinx Zynq-7000 SoC ZC706          @ 143 MHz

- CPU: Intel Comet Lake processor, 12 cores @ 2.9 GHz

- TX1: quad-core Arm Cortex-A57 processor @ 1.9 GHz

# Evaluation

*- Processing Latency and Energy of FPGA, CPU, and GPU*



- FPGA: Xilinx Zynq-7000 SoC ZC706        @ 143 MHz

- CPU: Intel Comet Lake processor, 12 cores @ 2.9 GHz

- TX1: quad-core Arm Cortex-A57 processor @ 1.9 GHz

# Evaluation

*- Processing Latency and Energy of FPGA, CPU, and GPU*



- FPGA: Xilinx Zynq-7000 SoC ZC706      @ 143 MHz
- CPU: Intel Comet Lake processor, 12 cores @ 2.9 GHz
- TX1: quad-core Arm Cortex-A57 processor @ 1.9 GHz

# Evaluation

*- Processing Latency and Energy of FPGA, CPU, and GPU*



| EuRoC Dataset (For drone) | FPGA Speedup | | FPGA Energy Reduction | |
|---|---|---|---|---|
| | Over CPU | Over TX1 | Over CPU | Over TX1 |
| FPGA ZC706 | 8.73x | 70.10x | 164.40x | 40.84x |
| Kintex-7 Series (XC7K160tfbg484) | 7.01x | 56.30x | 180.73x | 44.90x |
| Virtix-7 Series (XC7VX690tffg1761) | 10.75x | 86.34x | 172.05x | 42.75x |

| KITTI Dataset (For car) | FPGA Speedup | | FPGA Energy Reduction | |
|---|---|---|---|---|
| | Over CPU | Over TX1 | Over CPU | Over TX1 |
| FPGA ZC706 | 10.49x | 45.48x | 182.88x | 24.51x |
| Kintex-7 Series (XC7K160tfbg484) | 8.27x | 35.82x | 196.09x | 26.28x |
| Virtix-7 Series (XC7VX690tffg1761) | 12.71x | 55.08x | 188.60x | 25.28x |

# Evaluation

*- Comparison with Related Work*

| | This work | ISSCC'19 CNN-SLAM [1] | JSSC'19 Navion [2] | TC'20 pi-BA [3] | RSS'17 VIO on Chip [4] | HPCA'21 Eudoxus [5] |
|---|---|---|---|---|---|---|
| Platform | FPGA | ASIC | ASIC | FPGA | FPGA | FPGA |
| Technology | 28 nm | 28 nm | 65 nm | 28nm | 28nm | 16nm |
| Design | digital | digital | digital | digital | digital | digital |
| Type | SLAM | SLAM | SLAM | SLAM | SLAM | SLAM |
| Algorithm | Levenberg-Marquardt (optimization-based) | Levenberg-Marquardt (optimization-based) | Gaussian-Newton (optimization-based) | Levenberg-Marquardt (optimization-based) | Gaussian-Newton (optimization-based) | Kalman Filter (Filter-based) |
| DoF | 6-DoF | 6-DoF | 6-DoF | 6-DoF | 6-DoF | 6-DoF |
| Voltage | 1 V | 0.63-0.9V | 1.2V | 1 V | 1 V | 0.85 V |
| Power | 3.45W | 243.6mW @ 0.9V 61.75mW @ 0.63V | 24mW | 5.50W | 1.46 W | 8.96W |
| Frequency | 143 MHz | 240 MHz | 62.5/83.3 MHz | 143 MHz | 100 MHz | 180 MHz |
| Throughput | 55.8 GOPS | 879.6 GOPS @ 0.9V 329.8 GOPS @ 0.63V | 10.5-59.1 GOPS | N/A | 4.4-24.6 GOPS | N/A |
| Latency | 16.43 ms | N/A | 30.8 ms | 110 ms | 200 ms | 44.6 ms |
| Energy per Frame | 56.6 mJ | N/A | 739.2 uJ | 605 mJ | 292 mJ | 399.6 mJ |
| Dynamic Optimiza-tion | Yes | N/A | N/A | No | No | No |

# Outline

- SLAM: Simultaneously Localization & Mapping
- Hardware Architecture
- Main Contributions
- Evaluations and Comparisons
- Summary

# Summary

- **Energy-efficient** and **runtime-reconfigurable** FPGA accelerator for robotic localization and mapping.

# Summary

- **Energy-efficient** and **runtime-reconfigurable** FPGA accelerator for robotic localization and mapping.

- Leverage data sparsity, locality, and parallelism inherent in localization.
  - **4.1x** memory reduction with symmetry and sparsity
  - **5.7x** compute time reduction with time-multiplexed and pipeline processing
  - **5.8x** power reduction with runtime reconfiguration and clock gating

# Summary

- **Energy-efficient** and **runtime-reconfigurable** FPGA accelerator for robotic localization and mapping.

- Leverage data sparsity, locality, and parallelism inherent in localization.
  - **4.1x** memory reduction with symmetry and sparsity
  - **5.7x** compute time reduction with time-multiplexed and pipeline processing
  - **5.8x** power reduction with runtime reconfiguration and clock gating

- Our design is **2 orders of magnitude** more energy efficient than CPU and GPU.

CICC

# Reference



[Wan, CICC 2022]

# Reference



[Wan, Synthesis Lectures on Comp Arch 2021]     [Wan, CICC 2022]     [Wan, Circuits and Systems Magazine 2021]