

CatapultとTLM高位検証を用いた RISC-V Processorでの システムデザインの加速

Siemens EDA

TSS Calypto Japan

Consultant Application Engineer

酒井 健一

アジェンダ

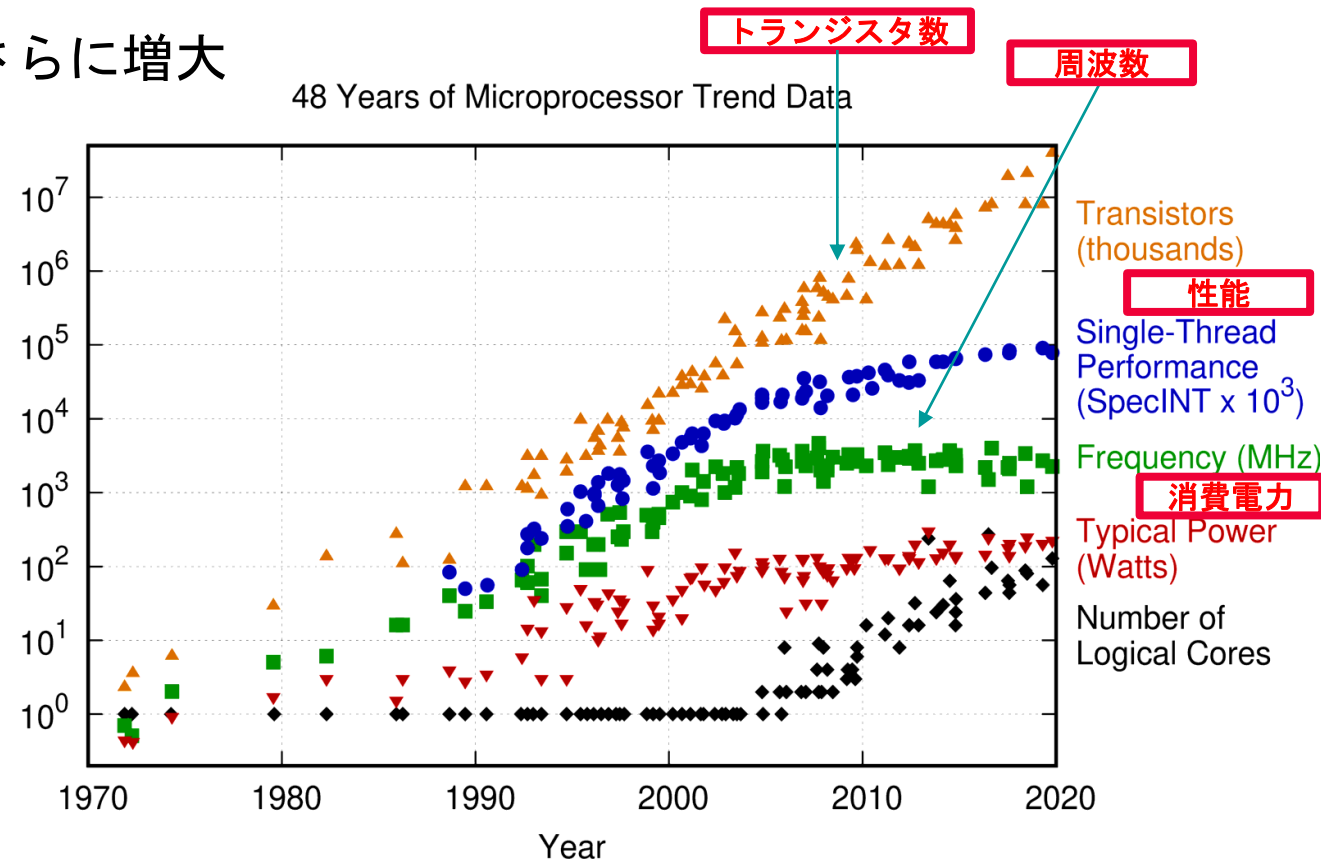
- プロセサを取り巻くトレンド
- プロセサモデルを含むシステムレベル検証の分類と目的
 - MatchLibの紹介
- Catapultが提供する最適解の探求とその事例
- まとめ

プロセッサを取り巻くトレンド

システム全体の最適化へのプロセッサSW要所の性能改善要求

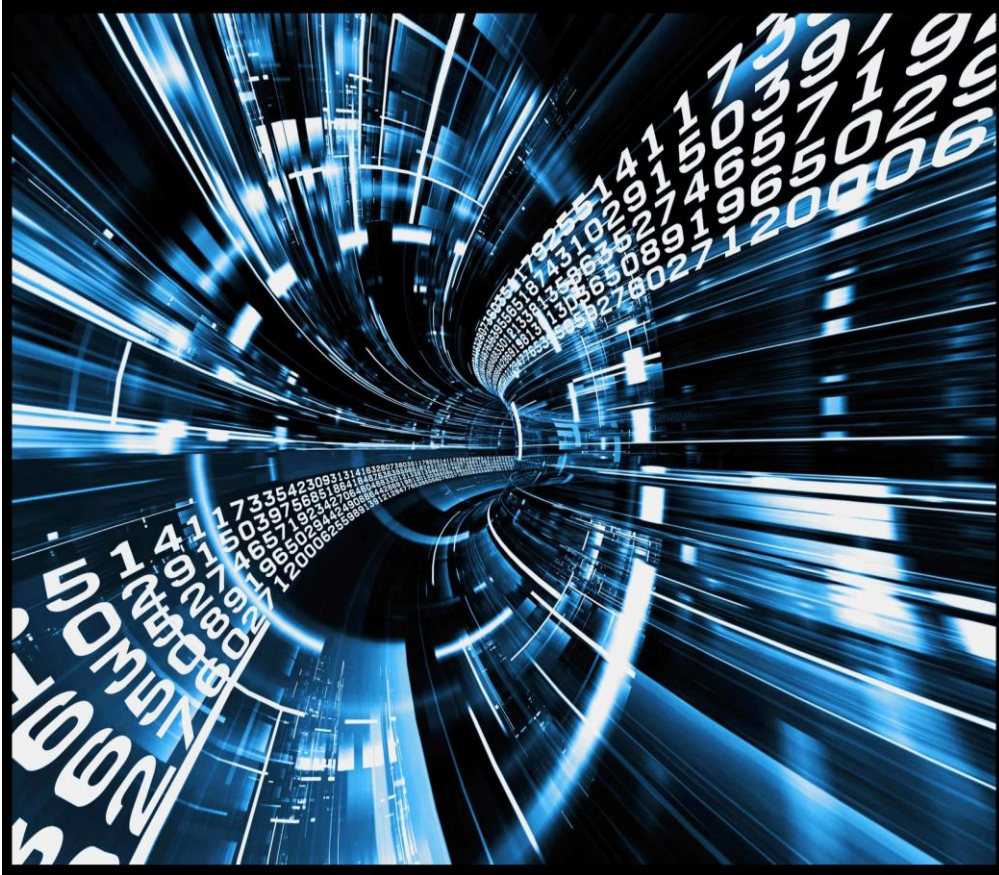
- ❑ ムーアの法則の終焉により、微細加工技術の進歩が鈍化
プロセッサの処理能力、エネルギー効率も鈍化
- ❑ 要求される処理能力、エネルギー効率は、さらに増大

このギャップを埋めるには？
理想論ではなく、現実的な解決策は？



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2019 by K. Rupp

HWを用いたアクセラレーション



SWはフレキシブル、プログラムが容易、迅速に更新

- 性能とエネルギーの面で非効率

HWへの機能の移動は、
性能とエネルギー面で著しい改善をもたらす

- 10X ~ 50Xの改善、 >80%のエネルギー削減

Video, 5G, AI アルゴリズムでは、
CPUが実行する以上にプロセッシングパワーを
必要とする

- これらの分野のアルゴリズムは、SWとして設計
- アルゴリズムは急な変更を伴い
RTL設計はリスクあり

HWへの機能のポーティング



HWが高速でより効率が良いか？

- プロセッサもHWで構成される
- CPUでの乗算機能 = HWの乗算機能

性能改善目的で、並列化を導入

- HWでは並列化の操作を実行することで、SW以上の速度を達成できる

効率のために、回路を調整する必要がある

- HWに必要なものだけを実装することでSWを上回る効率を提供できる

プロセッサモデルを含む システムレベル検証の 分類と目的

+ MatchLibの紹介

プロセサモデルを含むシステムレベル検証の分類と目的

プロセサモデル	インタコネク	Catapultの コンポーネント	検証の目的・タスク
Native C++	N/A	C++アルゴリズム or C++アーキテクチャ	機能面だけの検証 バスを介さない ポインタ制御でデータ移動
Host Code Execution (HCE)	TLM/MatchLib*	C++アーキテクチャ	インタフェースとバスの データ移動の考慮
Fast ISS	TLM/MatchLib*	C++アーキテクチャ	SWをターゲットCPUにクロスコンパイル 実装し、プロセサ性能を考慮したインタ フェースとバスのデータ移動の考慮 オープンソースや商用で利用可能 - > TLM 2.0への接続の鍵
RTL	RTL	高位合成後のRTL	クロック精度 最終製品での機能・性能検証

MatchLib : 次ページにて解説

MatchLibとは？

- **M**odular **A**pproach **T**o **C**ircuits and **H**ardware **L**ibrary
- MLアクセラレータを作成する際に、NVIDIA Labで開発
 - システムの挙動をシミュレーションするための**高抽象度**からのニーズ
 - (正確でなくとも)モデルの性能 ≒ RTL実装時の**性能** に極力近づけたいニーズ
- 性能と消費電力に対し、**多くの異なるアーキテクチャ**評価が必要
 - RTLで全てをデザインする余裕は無し
 - 重大な間違いが許容できる余裕は無し
- MatchLibは、オープンソースコード
 - <https://github.com/NVlabs/matchlib>

MatchLib

build passing

MatchLib is a SystemC/C++ library of commonly-used hardware functions and components that can be synthesized by most commercially-available HLS tools into RTL.

Doxygen-generated documentation can be found [here](#).

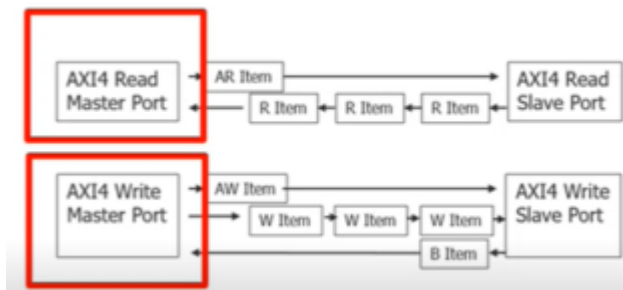
MatchLib is based on the Connections latency-insensitive channel implementation. Connections is included with the Catapult HLS tool and is available open-source on [HLSLibs](#). Additional documentation on the Connections latency-insensitive channel implementation can be found in the [Connections Guide](#).

合成可能なMatchLib AXI4 model と TLM Wrapper

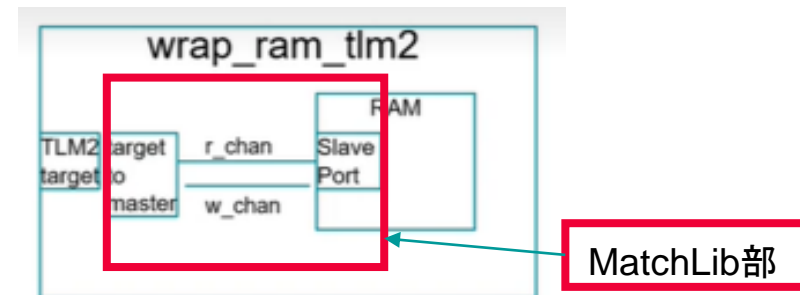
- AXI4の全5チャンネル (ar, r, aw, w, b) のtransactionを正確にモデル化
 - チャンネルでは、ビット精度とビート精度をモデル化
- 高いQoRを保ち高位合成Catapultで合成可能
- スループット精度 (ACCURATE_SIM)とFAST_SIMモードをサポート

- TLM < - > Catapult Wrapper Modelで、システム性能検証が早期から実現可能

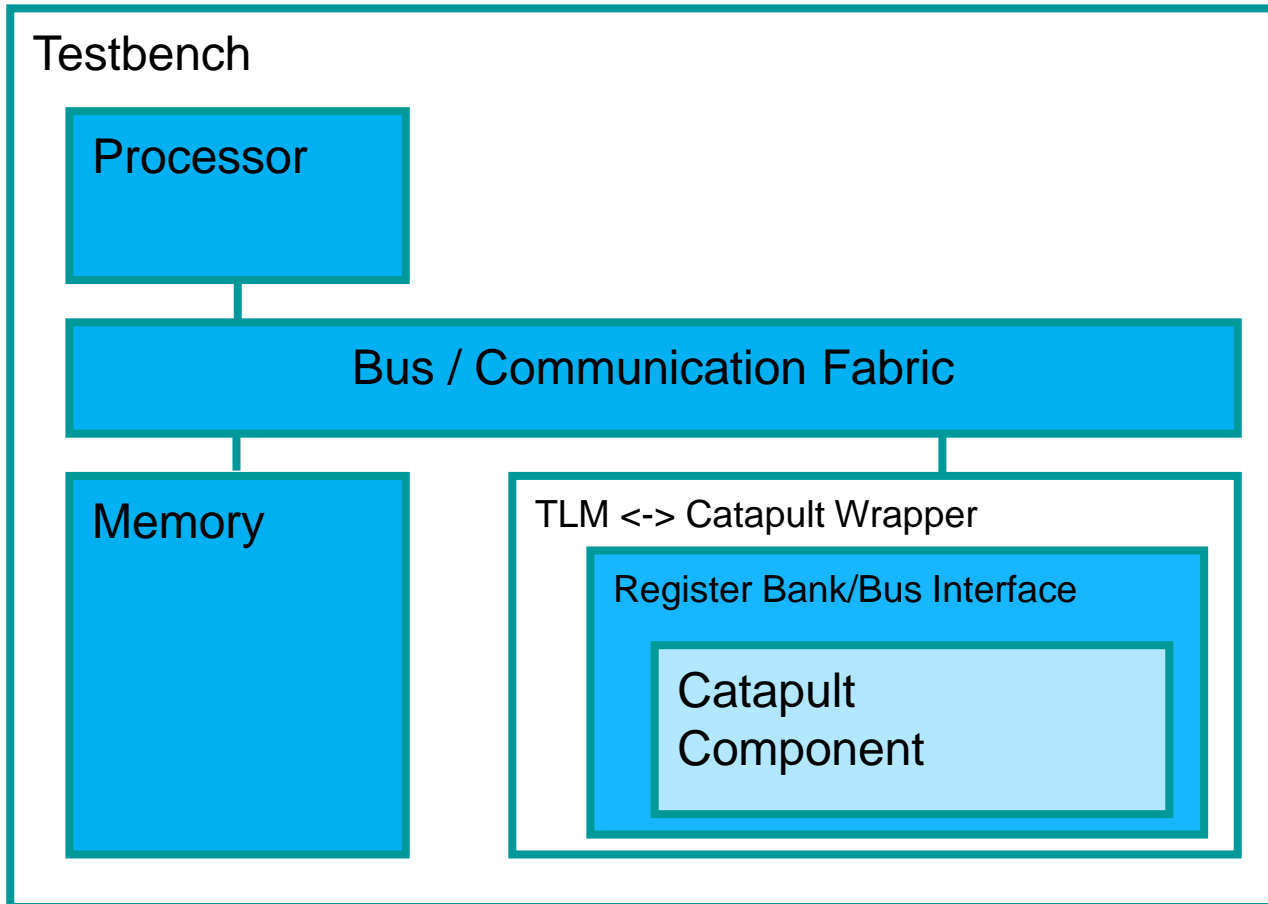
MatchLib Model



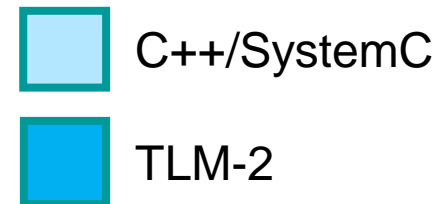
TLM Wrapper model for MatchLib Ram Model



Catapult向けのC++/SystemCのTLM環境への検証実装法



- ビット精度を扱うac_type から TLMへの変換を行うトランザクタで Catapultのコンポーネントをラッピング
- それ以外のシステムは TLM-2のコンポーネントで実装
 - コンカレントな実行動作



システムレベルシミュレーションの結果

(Yolo Tiny Inference)

検証目的	プロセサ	インタコネクト	アクセラレータ	RunTime (sec)	検証内容	速度	消費電力効率
機能検証 データと オペレータサイジング	Native C++	N/A	C++ Algorithmic	12	コアの アルゴリズム検証		
データ移動を伴う SWインタフェース/ 同期検証	Host Code	CONNECTIONS_FAST_SIM & TLM/MatchLib	C++ Architecture	1,027	インタフェース、 データフロー検証		
	Host Code	CONNECTIONS_ACCURATE_SIM & TLM/MatchLib	C++ Architecture	6,455			
ターゲットCPU SWインタフェース 性能検証 (approximate)	Fast ISS	CONNECTIONS_ACCURATE_SIM & TLM/MatchLib	C++ Architecture	7,272	クロスコンパイル コードでの SWドリブン検証	データ移動に依存 桁違いの性能差 システムの間 データ格納・アク セスに大きく依存	
最終性能と機能検証	RTL	CONNECTIONS_ACCURATE_SIM	C++ Architecture	97,329			
	RTL	RTL	RTL	(est) 2,600,000		クロック 精度	PowerPro +/- 15%
	Gate, post P&R	Gate, post P&R	Gate, post P&R	不明 アクセラレータが 必要			PowerPro +/- 2-5%

検証の目的に依るが、
SWドリブンの検証
では時間が掛かる

システムレベル検証からの考察

1. Host Code / ISS or QEMUモデル + MatchLib TLM高位検証により、ソフトウェア（プロセッサ処理）/データ移動のボトルネックを早期の時点で探ることが重要



2. 上記ボトルネックを探ることで、アクセラレーションの候補を見つけ出せる



3. アクセラレーションの手段として、高位合成によりアルゴリズムやアーキテクトの探求により、システムの最適解を見つけ出せる可能性が最大限に高まる

Catapultが提供する 最適解の探求とその事例

Catapult HLSが提供する最適解の探求のための機能

ユーザへのアーキテクチャ制御

- 並列化、スループット、面積、レイテンシ（ループの展開やパイプライン実装）
- メモリ vs レジスタ (リソース配置)

適用した制約でのアーキテクチャ探索と実装

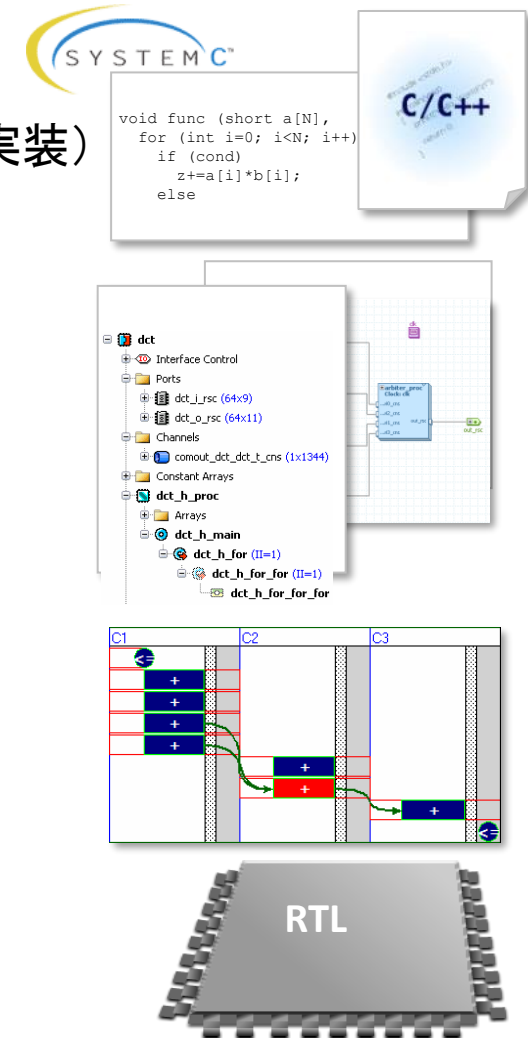
- ソースコードの変更を伴わない

自動化された算術最適化とビット幅最適化

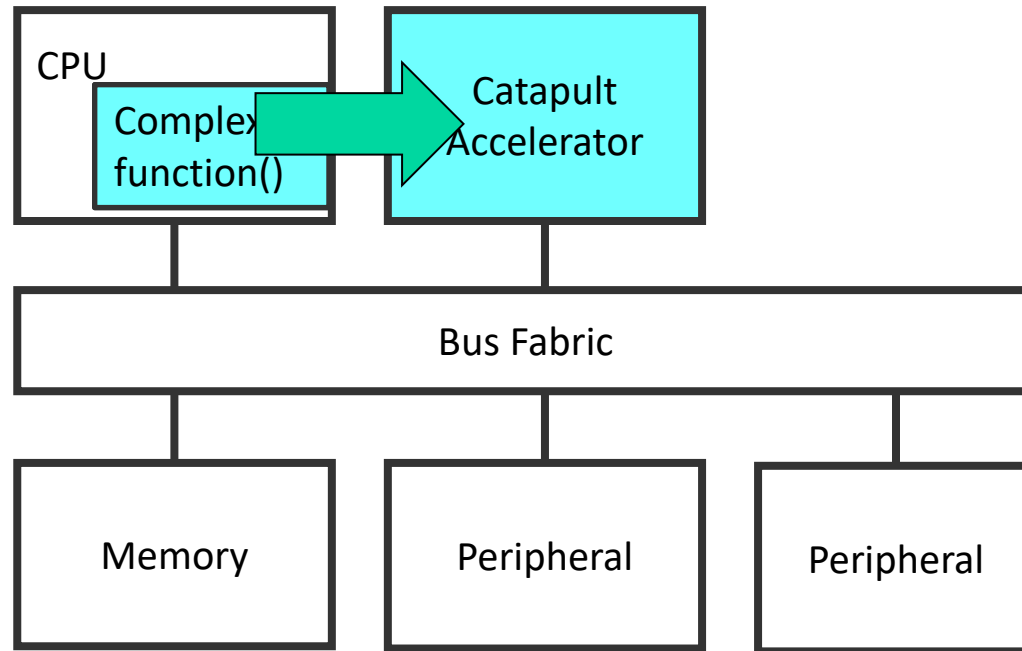
- ビット精度タイプにより、出力に算術精度の伝搬を可能に

FPGA とASICの双方にプロセスアウェアなスケジューリングを提供

- Area/Latency重視のデータパススケジューリングを実行
- I.P. 再利用のRTLテクノロジーのペナルティを排除



アクセラレーションの狙い



- CPUからHWに負荷のある機能の負荷を低減
- 性能の改善
- エネルギー消費の低減

プロファイルの取得と理解は大切！ アクセラレータ化前に...

Time Profiler > Profile > Root

Weight	Self Weight	Symbol Name
5.23 s 100.0%	0 s	▼yolo (48794)
5.23 s 100.0%	0 s	▼Main Thread 0x461821
5.23 s 100.0%	0 s	▼start libdyld.dylib
5.23 s 100.0%	0 s	▼main yolo
5.23 s 100.0%	0 s	▼yolo_tiny yolo
5.21 s 99.5%	4.99 s	►conv_2d_sw yolo
8.00 ms 0.1%	8.00 ms	leaky_relu_sw yolo
8.00 ms 0.1%	8.00 ms	max_pool_sw yolo
6.00 ms 0.1%	6.00 ms	bias_add_sw yolo

何処？
何が？

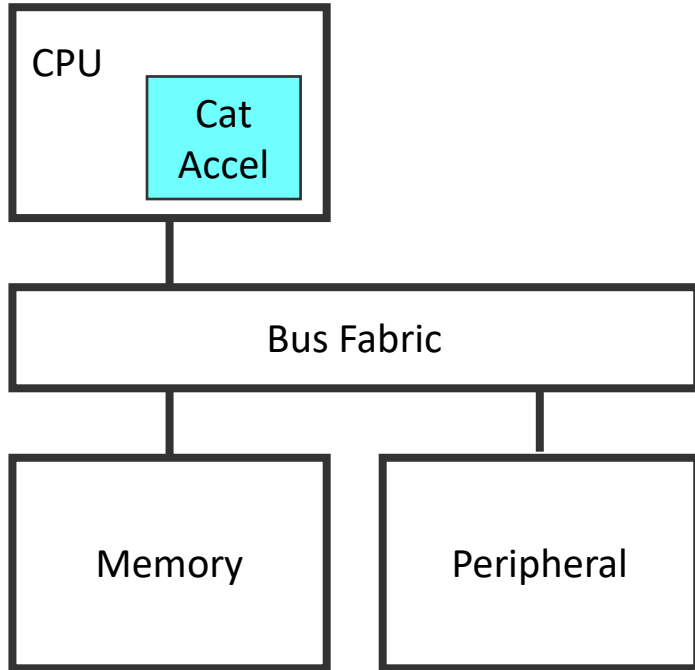


□ 計算上ボトルネックになる箇所を理解する

- 小規模だが頻繁な繰り返しのある関数は、
新規のインストラクションまたはco-processorに変更可能な候補
- 大規模な関数は、バスベースのペリフェラルに変更可能な候補
∴ 中間データの保持やデータ移動を伴うことになる



アクセラレータの接続例 – 新規のインストラクション



□ CPUと密にインテグレート

- バスファブリック上のトラフィックが無い
- 最速のオペランドアクセス

□ RISC-Vでは、インストラクションセットの拡張をサポート

□ 小規模だが複雑なオペレーション

- 数サイクルに限定
- 内部ステータを持たない

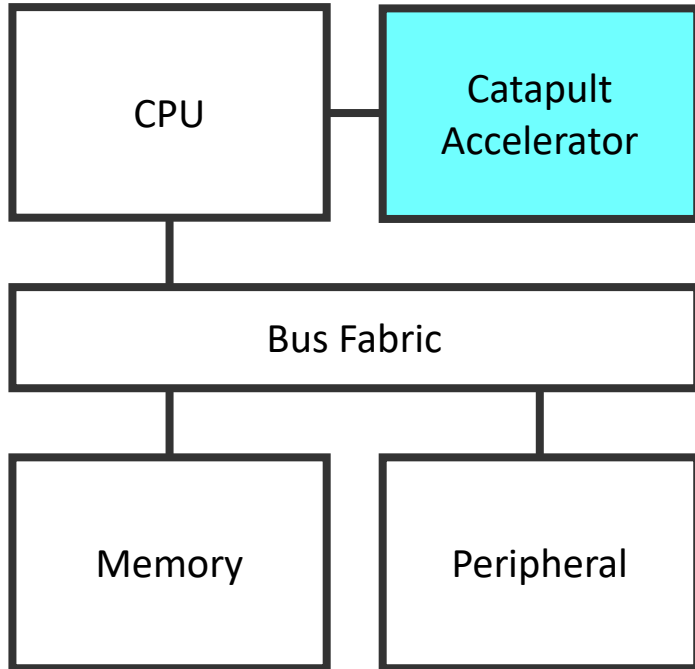
例:

- 複雑なSIMDインストラクション
- 複素数データタイプ演算

新規opcodeでのアクセスは、通常はC "asm"ディレクティブを使用

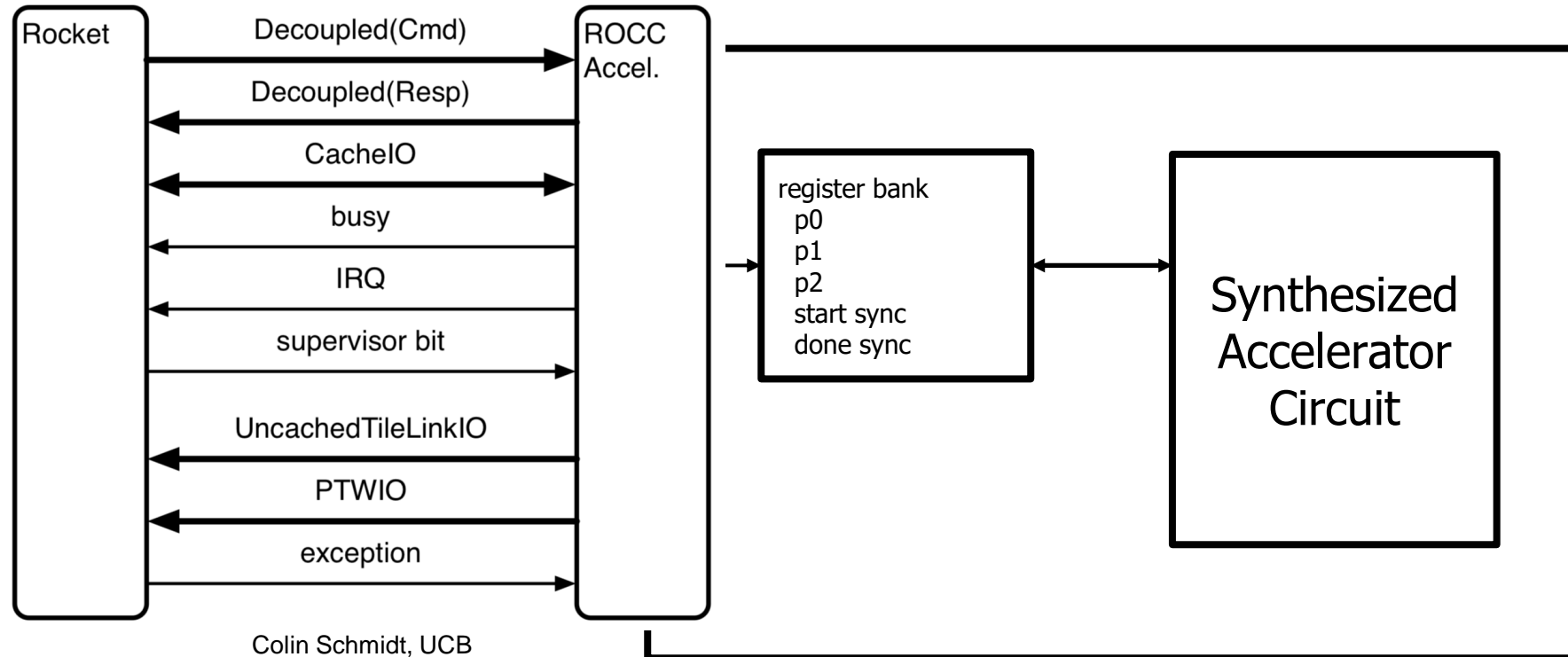
```
asm("CX3 %2, %0, %1"  
    : "r" var_in_0, var_in_1  
    : "=r" var_out_2)
```

アクセラレータの接続例– Co-processor (RoCC)

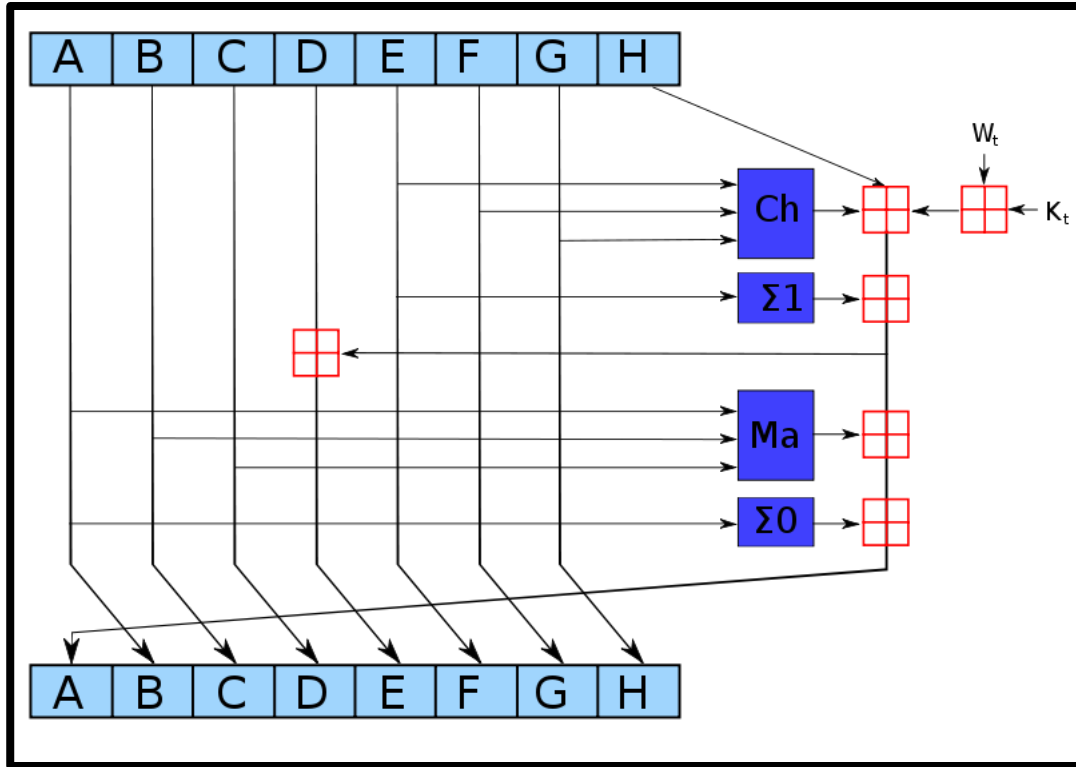


- CPUからco-processorインタフェースへデータを渡す
 - バスファブリック上のトラフィックが無し
 - 高速のアクセス
 - CPUだけがアクセラレータにアクセスできる。
- Rocket Core Custom Coprocessor Interfaceを使用
 - メモリの一貫性の問題を避けるため、メモリアクセスの実行は行わない
- 中程度の複雑度のアルゴリズムでは、SWと組み合わせたデータを使用
 - 例:
 - フィルタ
 - エンクリプション
 - ハッシュ

アクセラレータの接続例– Co-processor (RoCC)



SHA-256 Hash Algorithm



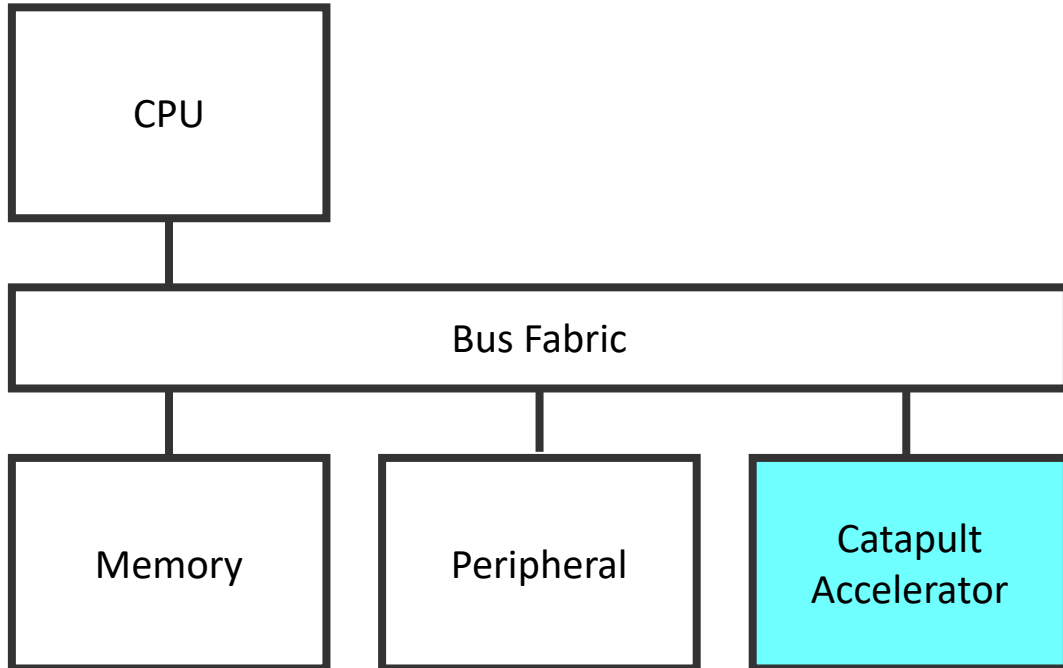
□ HW化により演算の複雑性を緩和

- 数100万のビット演算
- 高並列化
- データ移動は制限
- Hashのサイズ依存の性能

□ SWと比べたHW性能比

- 25X 高速
- 49X の低エネルギー (hash当たり)

アクセラレータの接続例- Slave

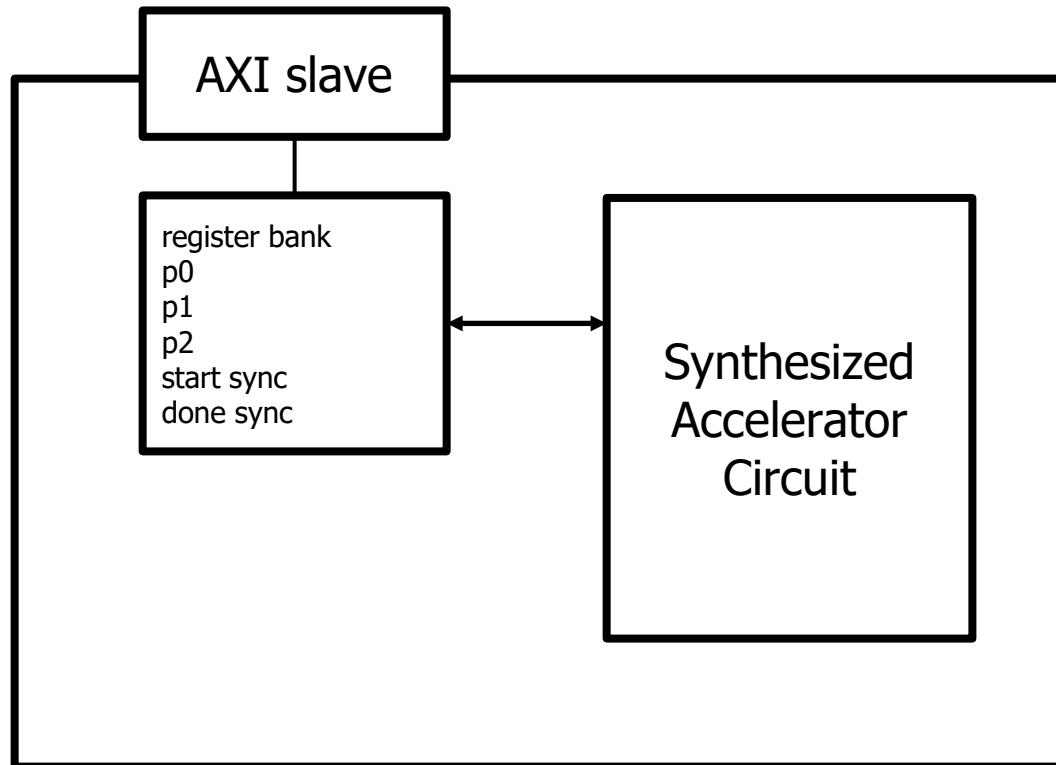


- データはCPUからアクセラレータへ渡される
 - ローカルメモリも実装可
 - キャッシュと仮想/物理アドレスの問題は無し
 - 低い複雑度のアルゴリズムでより高速
 - メモリマップドレジスタインタフェースを介してアクセス
- 低い複雑度の演算を軽減

例

- フィルタ
- エンクリプション
- ハッシュ

アクセラレータの接続例- Slave



複雑なAIアプリケーション 2D Convolution Algorithm

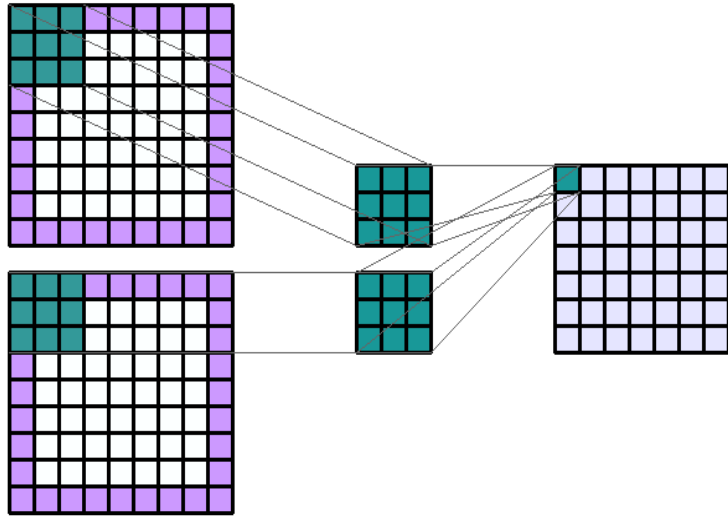
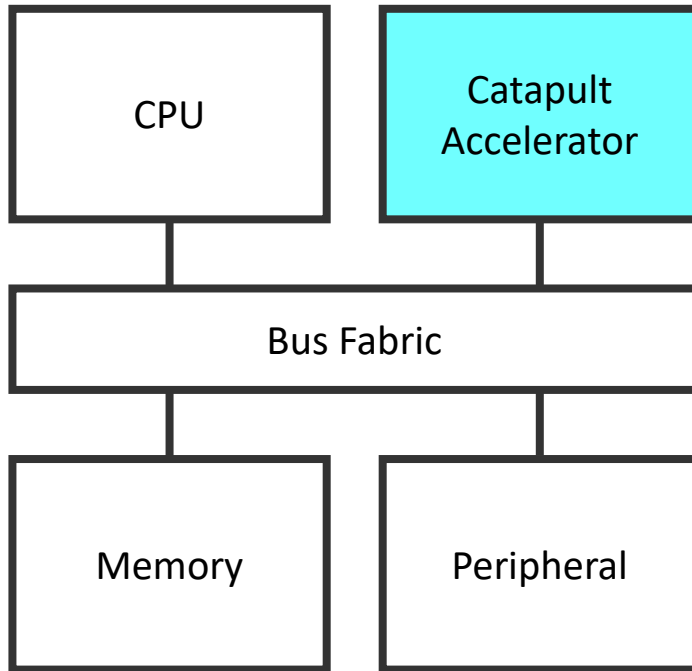


image source: towarddatascience.com

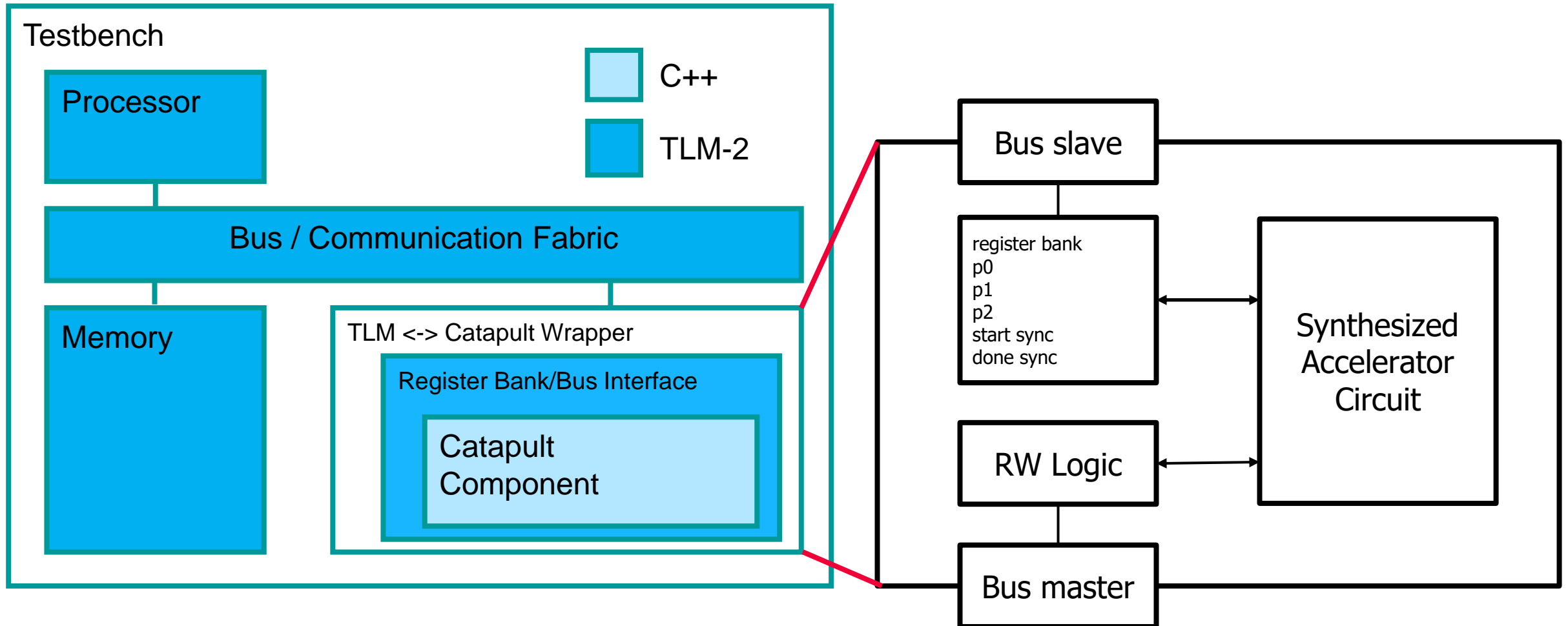
- 多くのAIと画像処理アプリケーションで使用
 - データ要素当たり数千の演算
 - 高並列性
 - データ量に依存する性能
 - 配列サイズとFeatureマップの数とフィルタの数に大きく依存
- SWと比べたHW性能比
 - 340X 高速
 - 480X の低エネルギー（演算当たり）

アクセラレータの接続例- Master

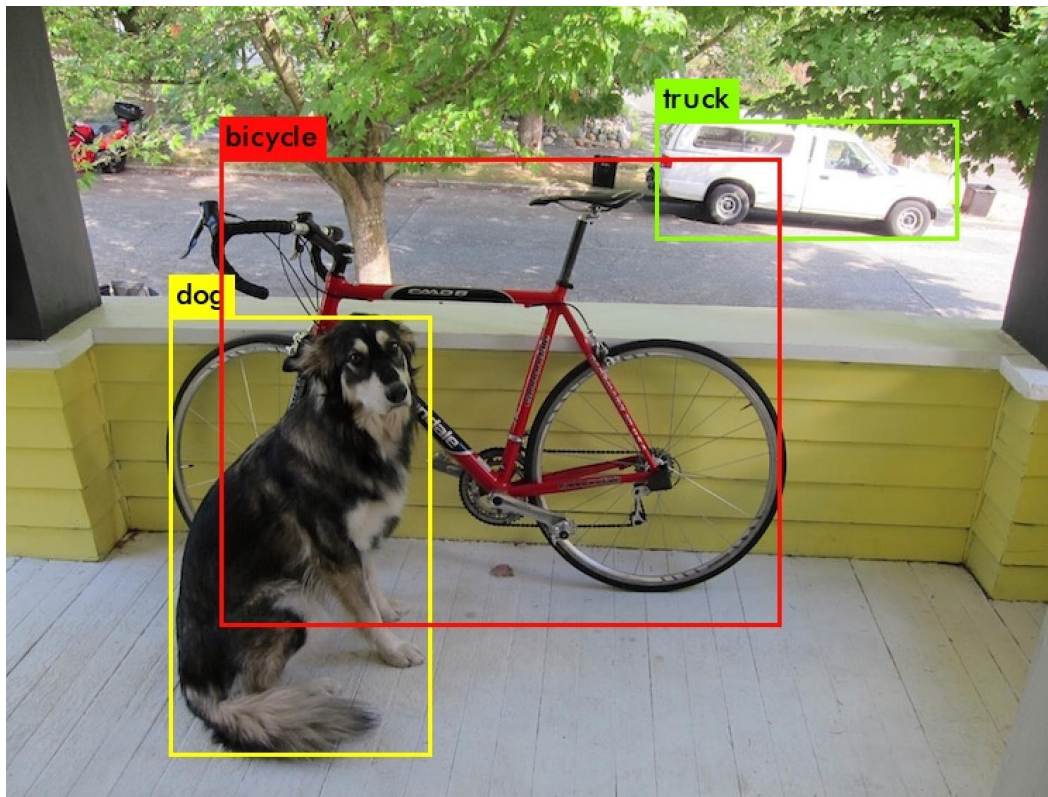


- アクセラレータがプロセサのデータメモリに直接アクセス可能
 - 入力の読み込み、出力の書き出し
 - 仮想/物理アドレスのキャッシュを考慮する必要あり
 - データの移動は性能を制限する
- コントロールと同期化のためのスレーブインタフェースを実装
 - メモリマップドレジスタインタフェースを介してアクセスされる
- 大規模で複雑な演算の負荷を低減する
 - 例:
 - 推論アルゴリズム
 - 画像 or ビデオプロセッシング
 - パケットフィルタリング

アクセラレータの接続例 - Master (+ Slave)



Yolo Tiny 物体認識アルゴリズム



□ 高程度の演算の複雑性

- 推論に~34億の MAC
- 高並列性
- 大きな中間データセット

□ SWとの性能比較

- 1,123X 高速
- 1,480Xの低エネルギー（推論当たり）

各実装方法での性能・エネルギー比較

Algorithm	Interface	SW perf/ HW perf	SW energy/ HW energy
bfloat	Instructions	N/A	N/A
SHA-256 hash	Co-Processor	25	49
conv2d	Bus slave	340	480
Yolo Tiny	Bus master	1,123	1,480

高位検証での検討と
高位合成によるHW
化の効果

Performanceの数値では
プロセサと同期のデータ移動のon/offのオーバーヘッドを含む

- 大規模の関数で、大きな性能の改善が得られる
 - 並列化が必須 シリアルのアプローチでは効果が出ない

まとめ – Catapultを用いたシステムデザインの加速 –

- プロセサとSWの影響を理解するため、Host Code / ISS or QEMUモデルを用いた高位検証が重要
- プロセサモデルに、CatapultのMatchLib + TLMコンポーネントを組み合わせることで高位合成前後の双方で、機能的な検証、性能、面積、電力をシステムレベルの全体で検証が可能
 - 検証のゴールに合う適切なレベルの抽象度の選択が鍵
- Catapult HLS は、システムデザインを加速するための必要不可欠なソリューションを提供
 - SWからHWへアクセラレーションでは、Catapultがベストなソリューション



Early Birdウェビナーシリーズ - 高位設計ソリューション シーメンスEDAジャパン https://www.mentorg.co.jp/events/early_bird/

SIEMENS 業界別ソリューション ソフトウェア & 製品 トレーニング & サポート

Home > イベント > Early Birdウェビナーシリーズ - 高位設計ソリューション

Early Birdウェビナーシリーズ

Catapult高位合成および高位検証プラットフォームやPowerPro RTLローパワー・ソリューション、Precision FPGA論理合成ソリューションなどを中心に、一歩先の未来を切り拓くユニークな高位設計ソリューションについて多角的に解説するウェビナーシリーズ

イベント

ご関心のある方は、是非こちらのページもご参照ください。

CatapultやMatchLibに関する様々な情報が満載です。

Early Bird

概要

「Early Bird」ウェビナーシリーズでは、Catapult高位合成および高位検証プラットフォームやPowerPro RTLローパワー・ソリューション、Precision FPGA論理合成ソリューションなどを中心に、一歩先の未来を切り拓くユニークなEDAソリューションについて多角的に解説いたします。また、シーメンスEDA製品の紹介にとどまらず、5G、AI、IoT時代において、開発者の皆さまが新たな視点やインサイトを得るきっかけとなるような新しいテクノロジーや業界動向について、社外講師を迎えての講演も積極的にお届けします。

| Contact

Siemens EDA

TSS Calypto Japan

Consultant Application Engineer

酒井 健一

E-mail: Kenichi_sakai@mentor.com