

Qualcomm

May 9, 2023

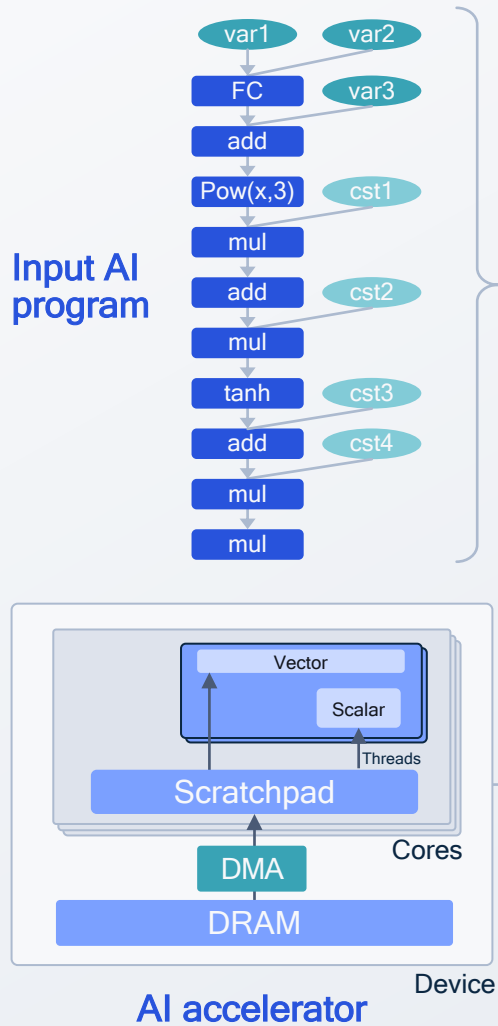
Advanced compilation technology for the AI stack

Benoît Meister

Principal Engineer
Qualcomm Technologies, Inc.

@QCOMResearch

Snapdragon and Qualcomm branded products are products of Qualcomm Technologies, Inc. and/or its subsidiaries.



Qualcomm®
Polyhedral
Mapper

Exploits key
hardware
speedup
mechanisms

Optimized code (sample)

```
void subgraph_3_1(float (*__restrict__ var)[1024],
float (*__restrict__ var_2)[1024],
float* __restrict__ var_3,
float (*__restrict__ T_multiply)[4096])
{
/* [var decls] */
_t1_1 = rspmd_core_id(0);
_t2 = rspmd_thd_id(0);
var_1 = (float (*)[1024])static_alloc(spad, 0, 52ul, 131072ul) /* var_1 */;
/* [more static allocs] */
if (_t1_1 <= 11) {
if (_t2 <= 0) {
int i;
for (i = 0; i <= 1; i++) {
int j;
int j_1;
int _t3;
int _t4;
rspmd_dma_get((void const volatile*)(var[32 * _t1_1] + 0), (void
volatile*)(var_1[32 * _t1_1 + -32 * _t1_1] + 0), (unsigned long)
(unsigned int)4096, (long)4096, (long)4096, (unsigned long)
(unsigned int)32, 0);
rspmd_dma_get((void const volatile*)(var_2[_2048 * i] + 0), (void
volatile*)(var_2_1[-2048 * i + 2048 * i] + 0), (unsigned long)
(unsigned int)4096, (long)4096, (long)4096, (unsigned long)
(unsigned int)2048, 0);
rspmd_dma_get((void const volatile*)(var_3_ + 2048 * i), (void
volatile*)(var_3_1 + (-2048 * i + 2048 * i)), (unsigned long)
(unsigned int)8192, (long)0, (long)0, (unsigned long)
(unsigned int)1, 0);
rspmd_dma_wait(0);
rspmd_thd_barrier();

```

```
for (j = 0; j <= 15; j++) {
int k;
for (k = 0; k <= 127; k++) {
int i1;
rv_vstore_f32t16(T_dense_comp1[_t2] + 0, rv_bcast_f32t16(0.0f));
for (i1 = 0; i1 <= 1023; i1++) {
rv_vstore_f32t16(T_dense_comp1[_t2] + 0, rv_fadd_f32t16(
rv_vload_f32t16(T_dense_comp1[_t2] + 0), rv_fmuls_f32t16(
(rv_vload_s_f32t16(var_1[j + 16 * _t2] + i1, 0),
rv_vload_s_f32t16(var_2_1[16 * k + 0] + i1, 1024))));
}
rv_vstore_f32t16(T_add_comp0 + 0, rv_fadd_f32t16(
rv_vload_f32t16(T_dense_comp1[_t2] + 0),
rv_vload_f32t16(var_3_1 + (16 * k + 0)));
rv_vstore_f32t16(T_power_comp0 + 0, rv_fpow_f32t16(
rv_vload_f32t16(T_add_comp0 + 0),
rv_bcast_f32t16(3.0f));
rv_vstore_f32t16(T_multiply_2_comp0 + 0, rv_fmuls_f32t16(
rv_bcast_f32t16(0.04471499845f),
rv_vload_f32t16(T_power_comp0 + 0)));
rv_vstore_f32t16(T_add_2_comp0 + 0, rv_fadd_f32t16(
rv_vload_f32t16(T_add_comp0 + 0),
rv_vload_f32t16(T_multiply_2_comp0 + 0));
rv_vstore_f32t16(T_multiply_3_comp0 + 0, rv_fmuls_f32t16(
rv_bcast_f32t16(0.7978850007f),
rv_vload_f32t16(T_add_2_comp0 + 0));
rv_vstore_f32t16(T_tanh_comp0 + 0, rv_ftanh_f32t16(
rv_vload_f32t16(T_multiply_3_comp0 + 0));
rv_vstore_f32t16(T_add_3_comp0 + 0, rv_fadd_f32t16(
rv_bcast_f32t16(1.0f),
rv_vload_f32t16(T_tanh_comp0 + 0));
rv_vstore_f32t16(T_multiply_4_comp0 + 0, rv_fmuls_f32t16(
rv_bcast_f32t16(0.5f),
rv_vload_f32t16(T_add_3_comp0 + 0));
rv_vstore_f32t16(T_multiply_1[j + 16 * _t2] + (16 * k + 0),
rv_fmuls_f32t16(rv_vload_f32t16(T_add_comp0 + 0),
rv_vload_f32t16(T_multiply_4_comp0 + 0));
}
}
rspmd_thd_barrier();
/* [...] */

```

Features: **Core/thread parallelism** Explicit data transfers (DMA) SIMD parallelism Other

Taking advantage of parallel processing for power efficiency
Enabling automatic optimization of AI code

Today's agenda

A woman with dark hair pulled back, wearing a white patterned blouse and light-colored trousers, is sitting in a dark leather chair. She is looking out a window at a city skyline at night, holding a smartphone in her hands. The scene is dimly lit, with light coming from the window, creating a contemplative atmosphere.

The need for improved AI compilers

Advancements with polyhedral compilation

Our leading compiler research

General directions

Q&A

Smartphone



Smart homes



Video conferencing



Autonomous vehicles



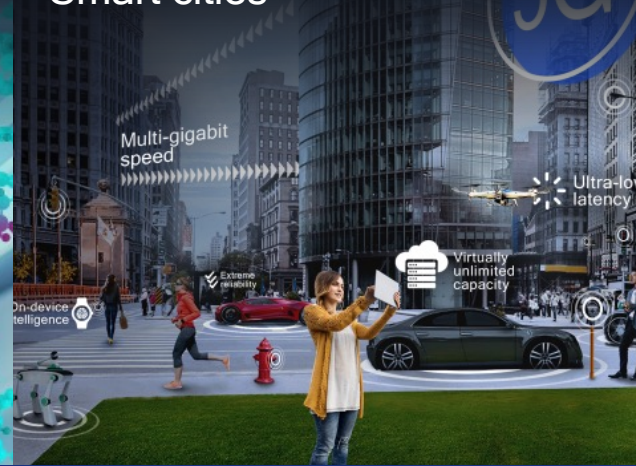
Smart factories



Extended reality



Smart cities



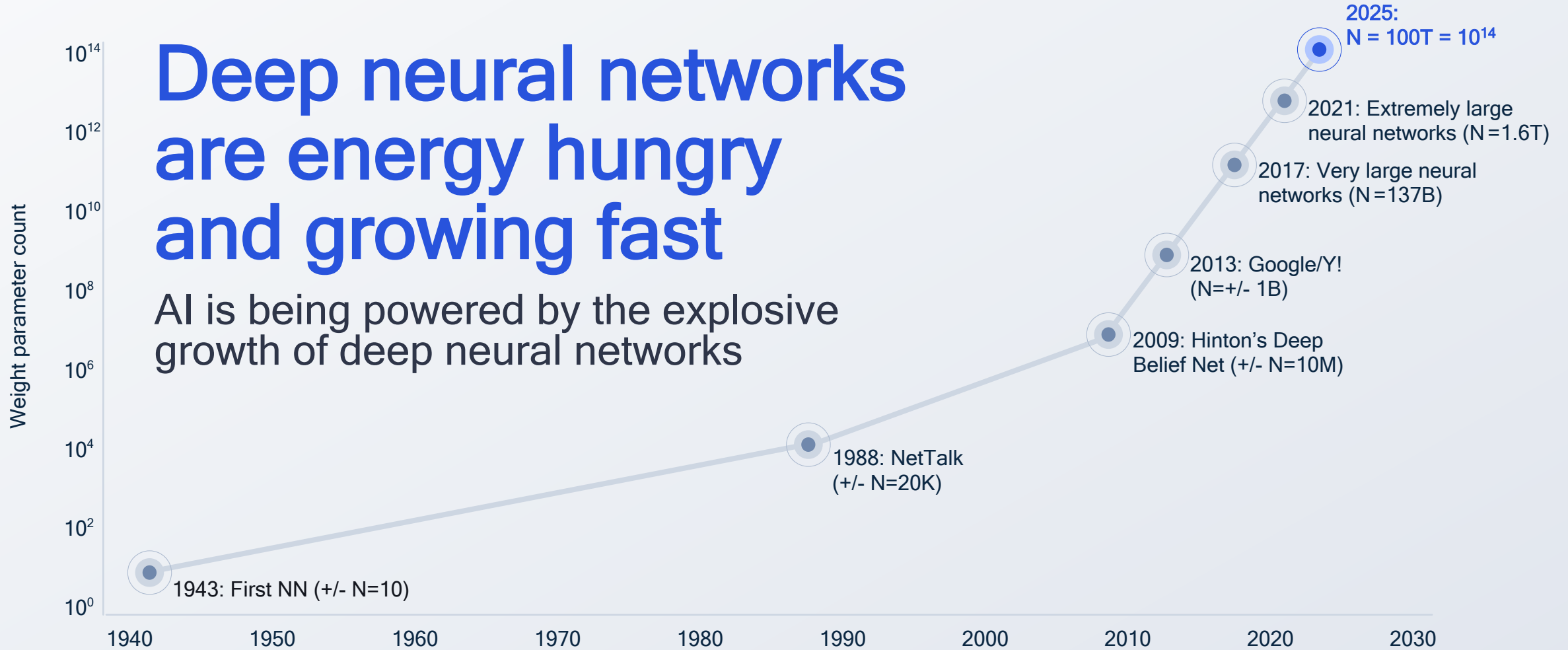
Video monitoring



AI is being used all around us
increasing productivity, enhancing collaboration, and transforming industries

Deep neural networks are energy hungry and growing fast

AI is being powered by the explosive growth of deep neural networks





2025


Will we have reached the capacity of the human brain?
Energy efficiency of the human brain is estimated to be 100,000x better than current hardware

The challenge of AI workloads

Constrained mobile environment

 Very compute intensive

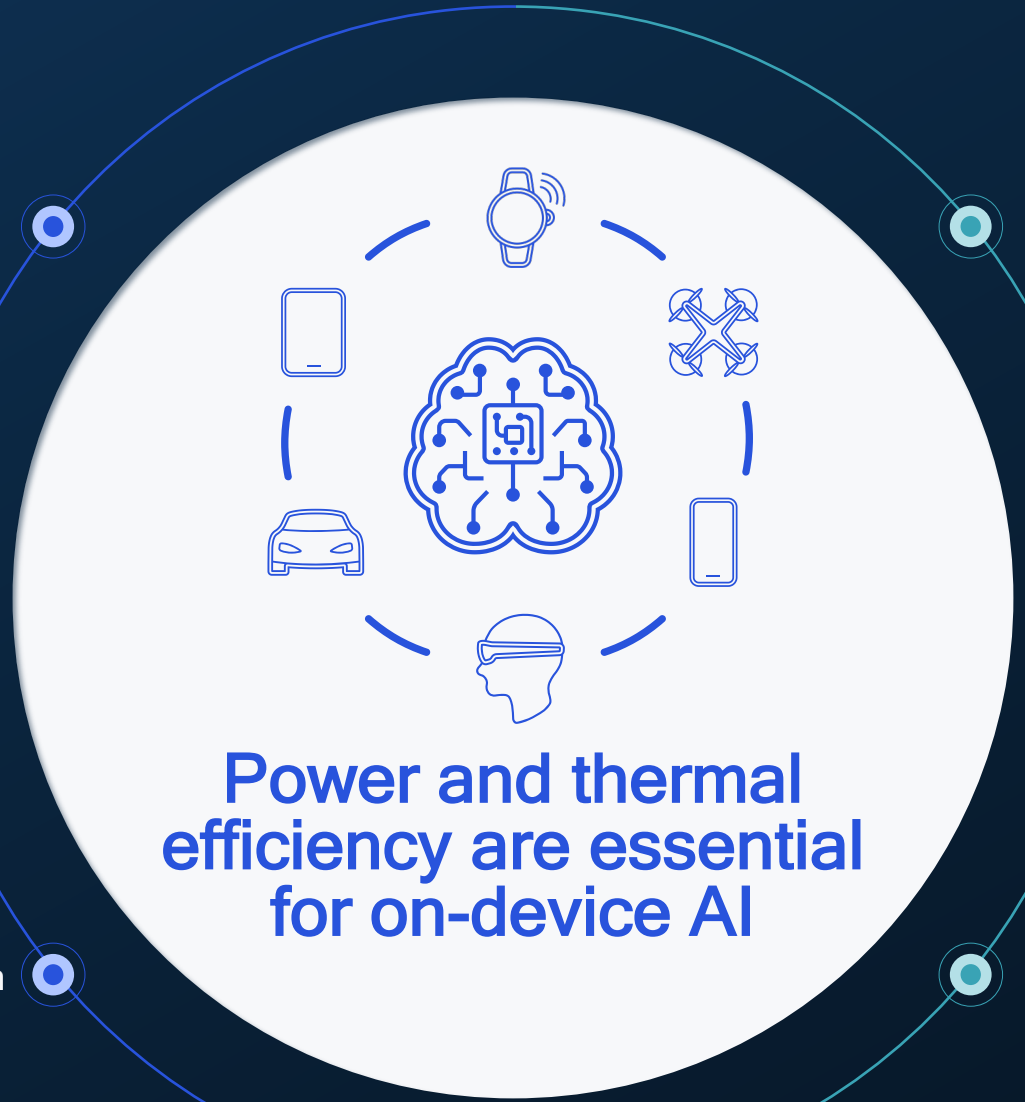
Must be thermally efficient for sleek, ultra-light designs 

 Large, complicated neural network models


 Complex concurrencies


 Real-time

 Always-on



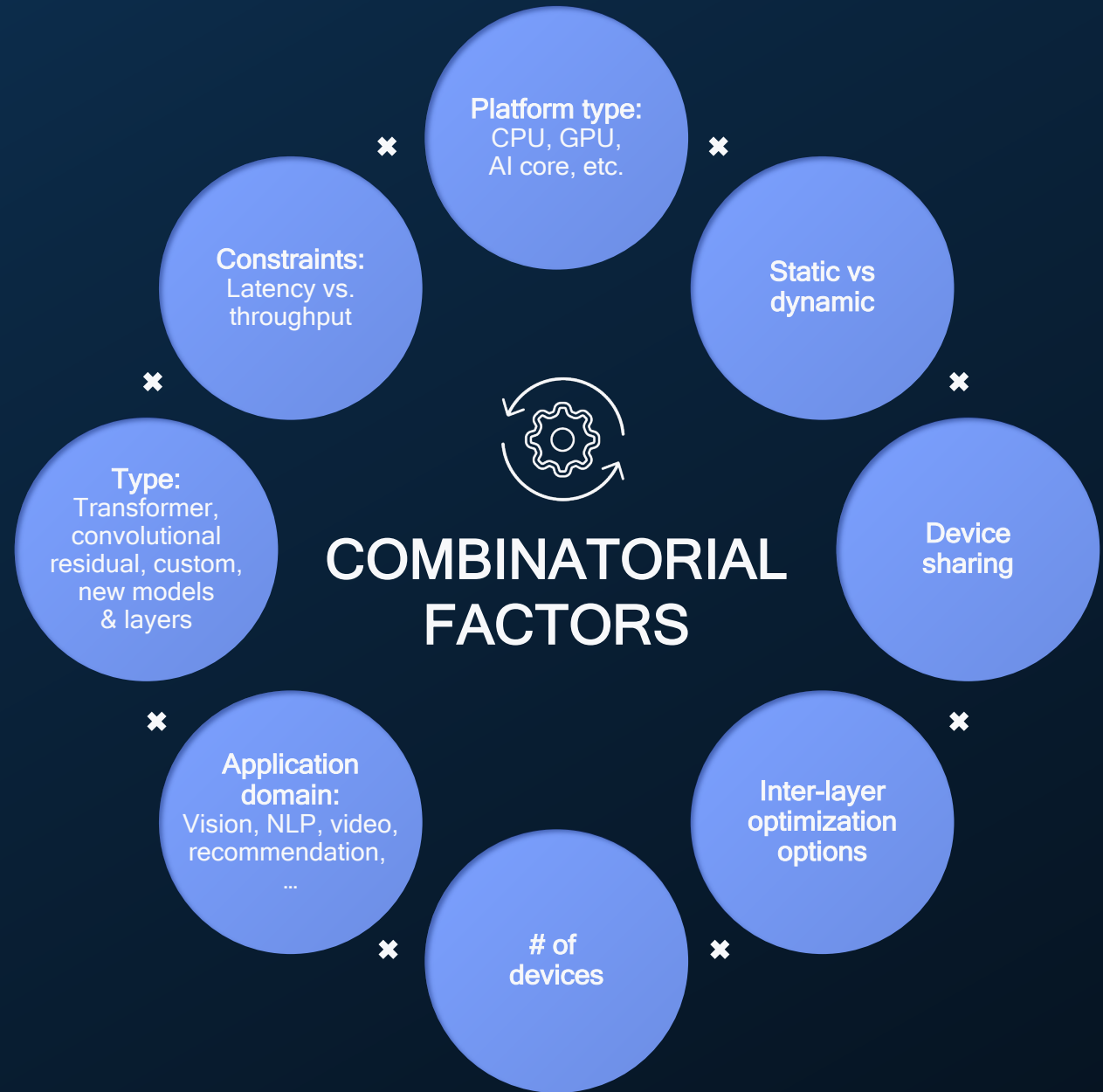
Power and thermal efficiency are essential for on-device AI

Requires long battery life for all-day use 

Storage/memory bandwidth limitations 

All optimization cannot be done manually

Optimization techniques need to scale with the explosion of AI



Holistic power efficiency research

Multiple axes to shrink AI models and efficiently run them on hardware

Quantization

Learning to reduce bit-precision while keeping desired accuracy

Conditional compute

Learning to execute only parts of a large inference model based on the input

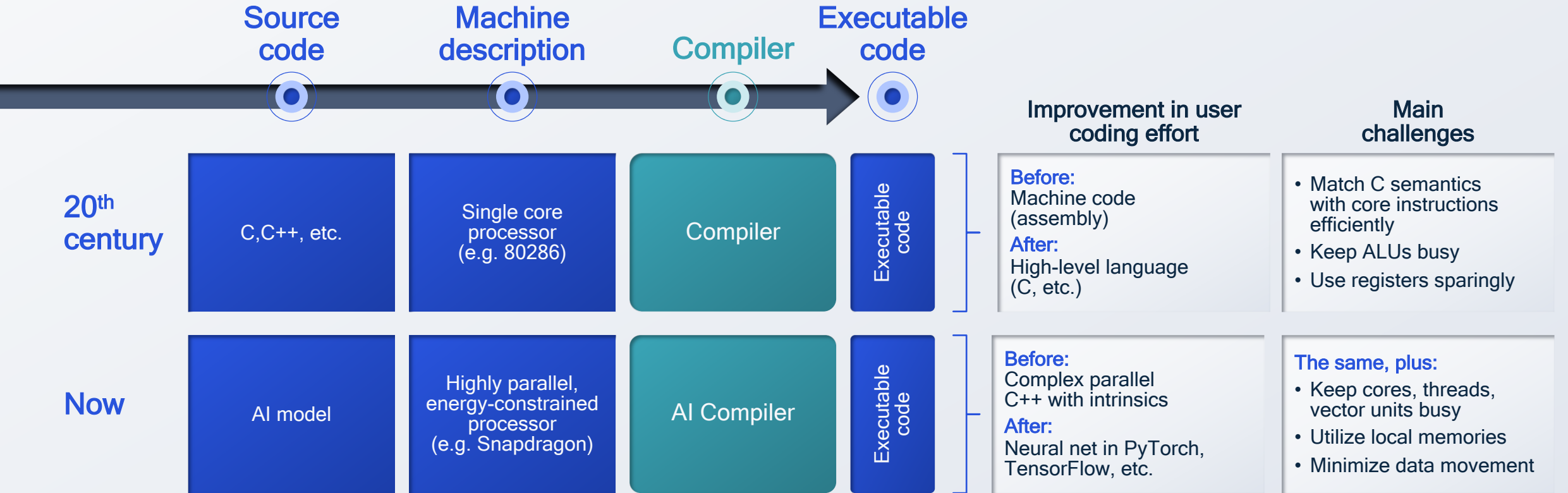
Neural architecture search

Learning to design smaller neural networks that are on par or outperform hand-designed architectures on real hardware

Compilation

Learning to compile AI models for efficient hardware execution

What is a compiler?



Compiler priorities have changed with evolving apps, software, and computer architectures

The need for parallelism and data locality

Use all the processing elements, avoid data transfers between remote and local memories

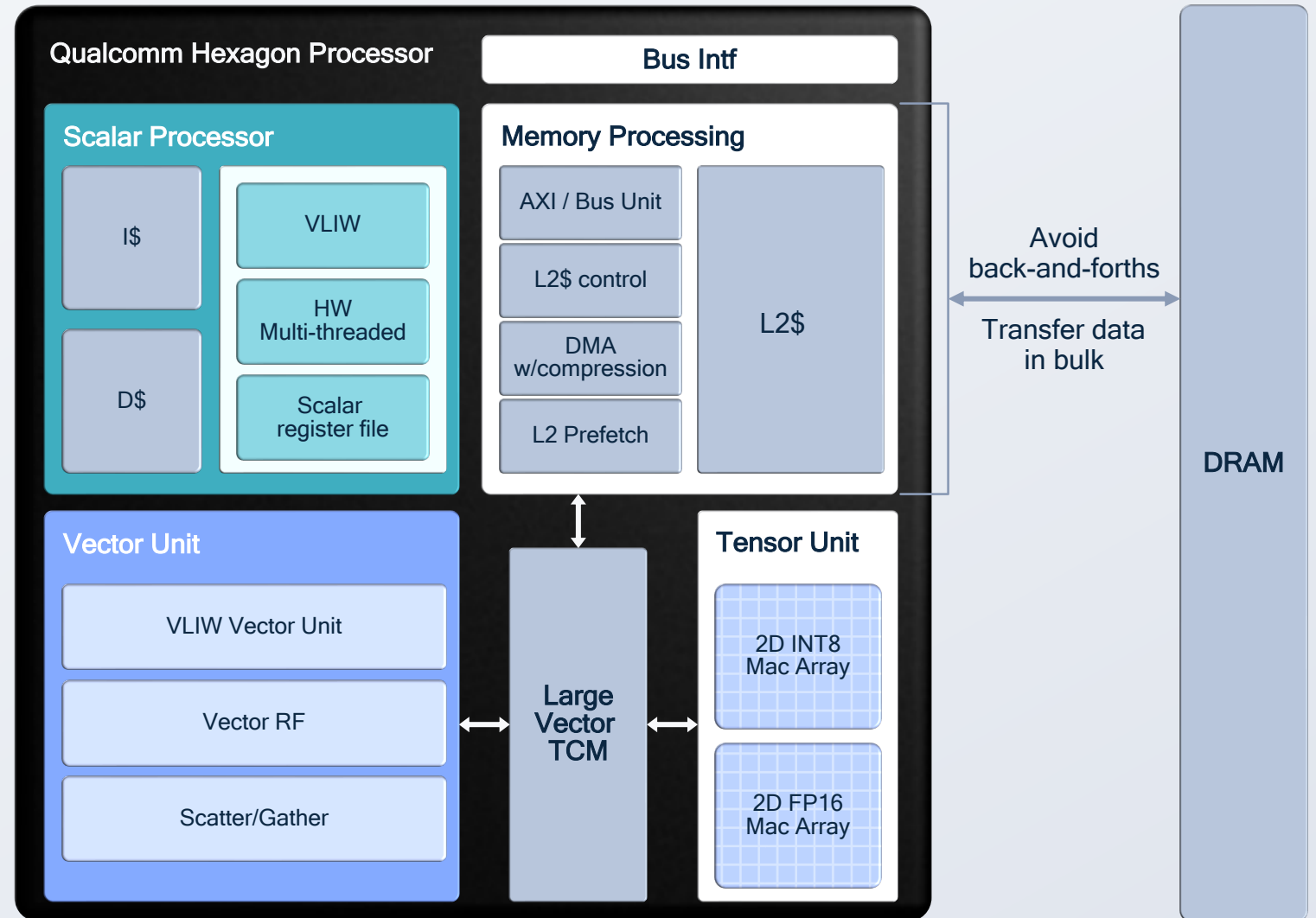
More PEs:
More compute, better power

Parallelism:
Exploit as many PEs as possible

Local memories:
Faster access to data

Data locality:
Use data shortly after it comes to local mem

Transfer useful data in bulk



Parallelism and data locality in code

Illustrated on a BERT language model subgraph optimized by the Qualcomm Polyhedral Mapper

```
void subgraph_3_1(float (*__restrict__ var)[1024],
float (*__restrict__ var_2)[1024],
float* __restrict__ var_3_,
float (*__restrict__ T_multiply)[4096])
{
    /* [var decls] */
    __t1_1 = rspmd_core_id(0);
    __t2 = rspmd_thd_id(0);
    var_1 = (float (*)[1024])static_alloc(spad, 0, 52ul, 131072ul) /* var_1 */;
    /* [more static allocs] */
    if (__t1_1 <= 11) {
        if (__t2 <= 0) {
            for (int i = 0; i <= 1; i++) {
                int j_1;
                int __t3;
                int __t4;
                rspmd_dma_get((void const volatile*)(var[32 * __t1_1] + 0), (void
volatile*)(var_1[32 * __t1_1 + -32 * __t1_1] + 0), (unsigned long)
(unsigned int)4096, (long)4096, (long)4096, (unsigned long)(
unsigned int)32, 0);
                rspmd_dma_get((void const volatile*)(var_2_[2048 * i] + 0), (void
volatile*)(var_2__1[-2048 * i + 2048 * i] + 0), (unsigned long)(
unsigned int)4096, (long)4096, (long)4096, (unsigned long)(
unsigned int)2048, 0);
                rspmd_dma_get((void const volatile*)(var_3_ + 2048 * i), (void
volatile*)(var_3__1 + (-2048 * i + 2048 * i)), (unsigned long)(
unsigned int)8192, (long)0, (long)0, (unsigned long)(
unsigned int)1, 0);
                rspmd_dma_wait(0);
                rspmd_thd_barrier();
            }
        }
    }
}
```

Features:

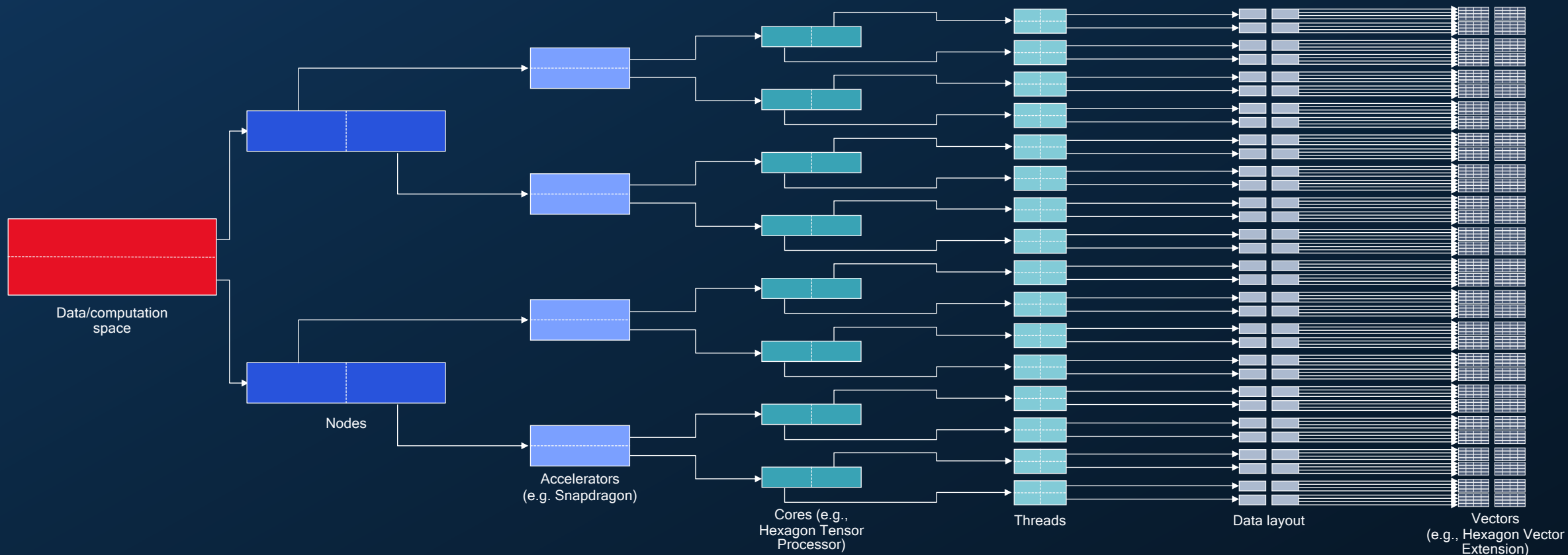
Core/thread parallelism Explicit data transfers (DMA) SIMD parallelism Other

```
for (int j = 0; j <= 15; j++) {
    for (int k = 0; k <= 127; k++) {
        rv_vstore_f32t16(T_dense_comp1[__t2] + 0, rv_bcast_f32t16(0.0f));
        for (int i1 = 0; i1 <= 1023; i1++) {
            rv_vstore_f32t16(T_dense_comp1[__t2] + 0, rv_fadd_f32t16(
rv_vload_f32t16(T_dense_comp1[__t2] + 0), rv_fmul_f32t16
(rv_vload_s_f32t16(var_1[j + 16 * __t2] + i1, 01),
rv_vload_s_f32t16(var_2__1[16 * k + 0] + i1, 10241))));
        }
        rv_vstore_f32t16(T_add_comp0 + 0, rv_fadd_f32t16(
rv_vload_f32t16(T_dense_comp1[__t2] + 0),
rv_vload_f32t16(var_3__1 + (16 * k + 0))));
        rv_vstore_f32t16(T_power_comp0 + 0, rv_fpow_f32t16(
rv_vload_f32t16(T_add_comp0 + 0),
rv_bcast_f32t16(3.0f)));
        rv_vstore_f32t16(T_multiply_2__comp0 + 0, rv_fmul_f32t16(
rv_bcast_f32t16(0.04471499845f),
rv_vload_f32t16(T_power_comp0 + 0)));
        rv_vstore_f32t16(T_add_2__comp0 + 0, rv_fadd_f32t16(
rv_vload_f32t16(T_add_comp0 + 0),
rv_vload_f32t16(T_multiply_2__comp0 + 0)));
        rv_vstore_f32t16(T_multiply_3__comp0 + 0, rv_fmul_f32t16(
rv_bcast_f32t16(0.7978850007f),
rv_vload_f32t16(T_add_2__comp0 + 0)));
        rv_vstore_f32t16(T_tanh_comp0 + 0, rv_ftanh_f32t16(
rv_vload_f32t16(T_multiply_3__comp0 + 0)));
        rv_vstore_f32t16(T_add_3__comp0 + 0, rv_fadd_f32t16(
rv_bcast_f32t16(1.0f),
rv_vload_f32t16(T_tanh_comp0 + 0)));
        rv_vstore_f32t16(T_multiply_4__comp0 + 0, rv_fmul_f32t16(
rv_bcast_f32t16(0.5f),
rv_vload_f32t16(T_add_3__comp0 + 0)));
        rv_vstore_f32t16(T_multiply_1[j + 16 * __t2] + (16 * k + 0),
rv_fmul_f32t16(rv_vload_f32t16(T_add_comp0 + 0),
rv_vload_f32t16(T_multiply_4__comp0 + 0)));
    }
}
rspmd_thd_barrier();
/* [ ... ] including __t2 > 0*/
}
```

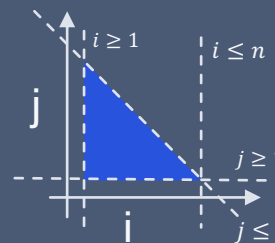
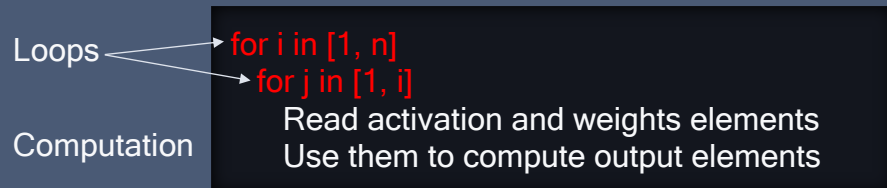
The Qualcomm Polyhedral Mapper

Qualcomm Technologies' polyhedral approach to AI model optimization

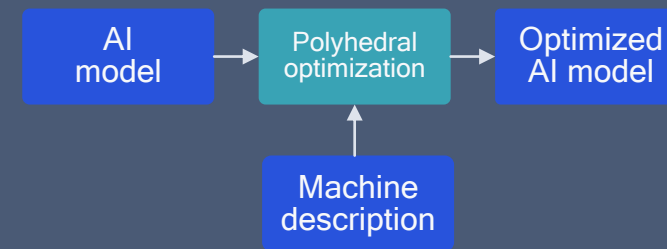
How we optimize AI models: the Qualcomm Polyhedral Mapper

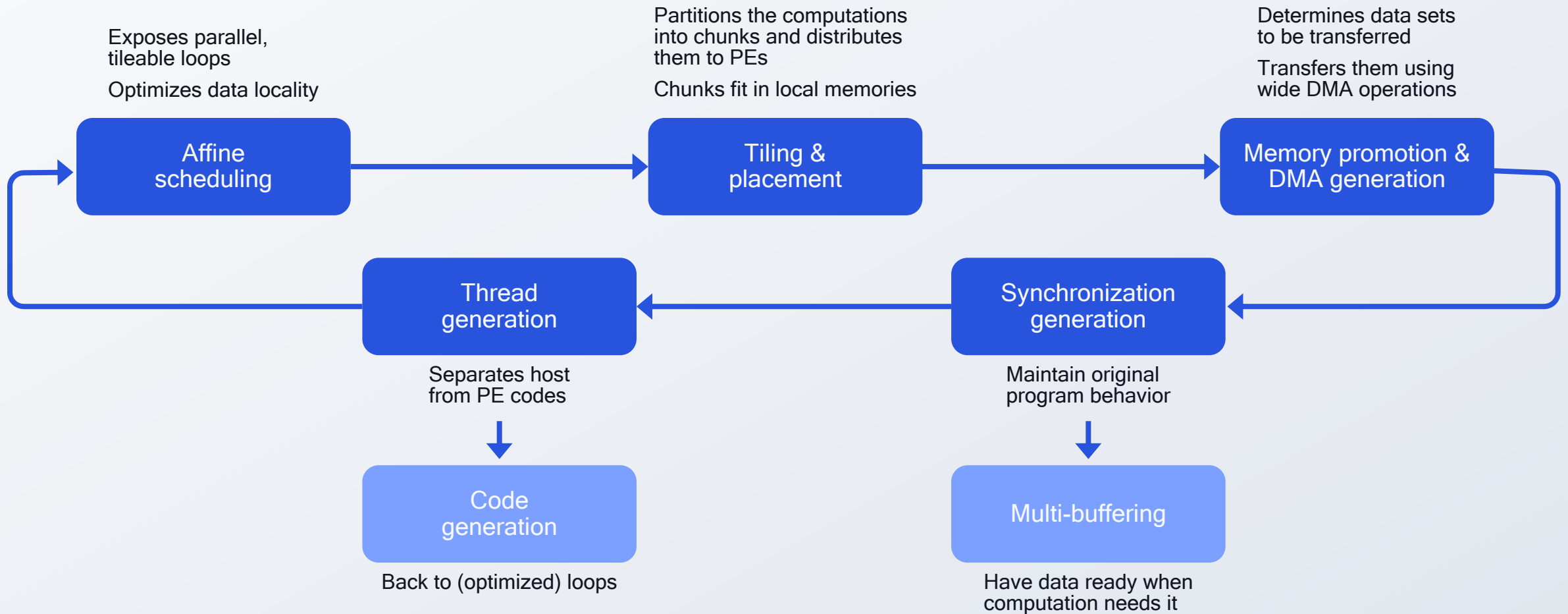


AI code is repeated computations



Model the loops as a polyhedron

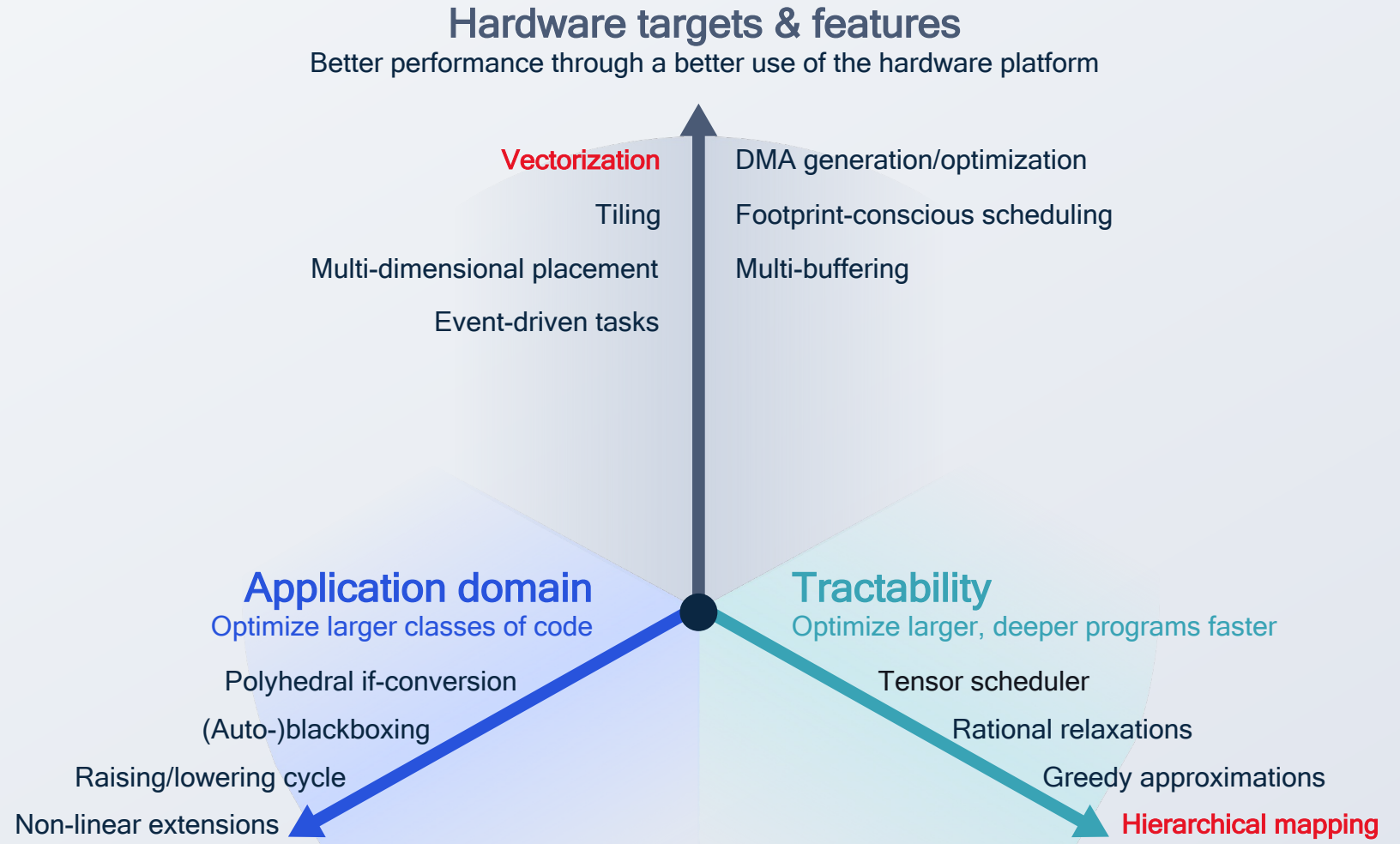




Bird's eye view of polyhedral AI compiler optimizations

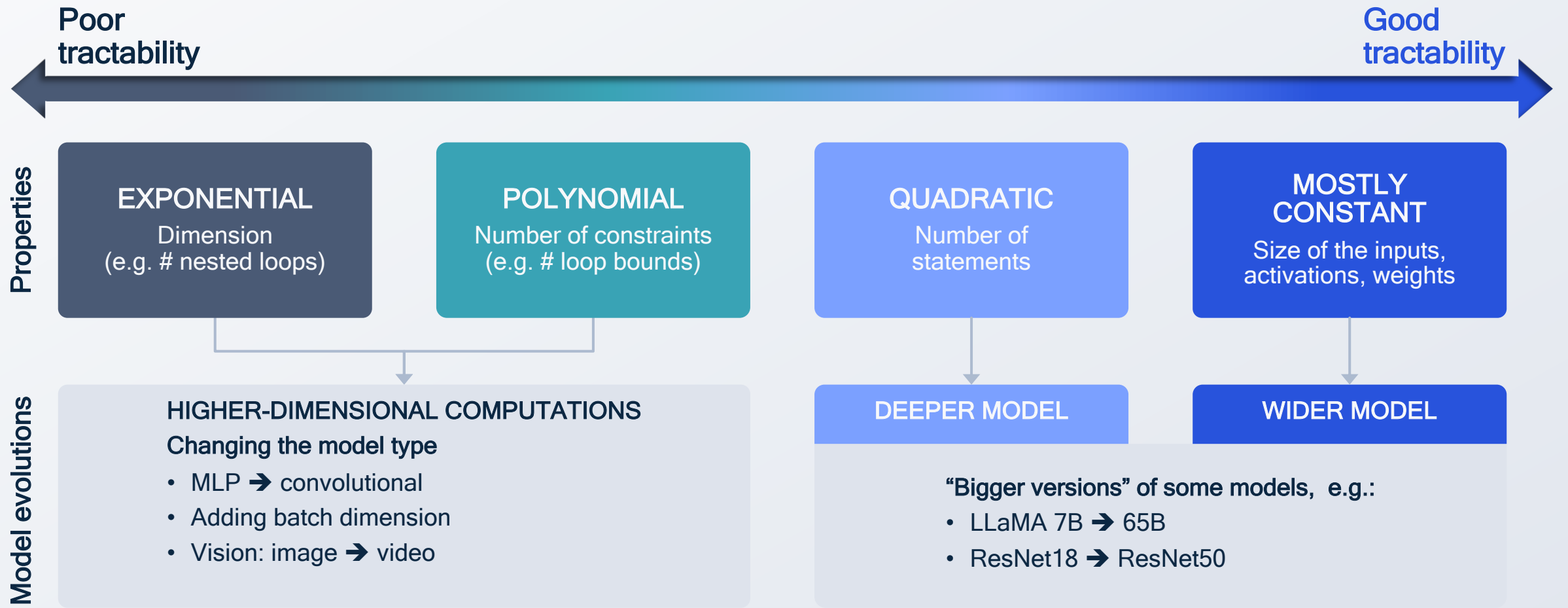
Challenge:

AI compilers need to grow in three different directions



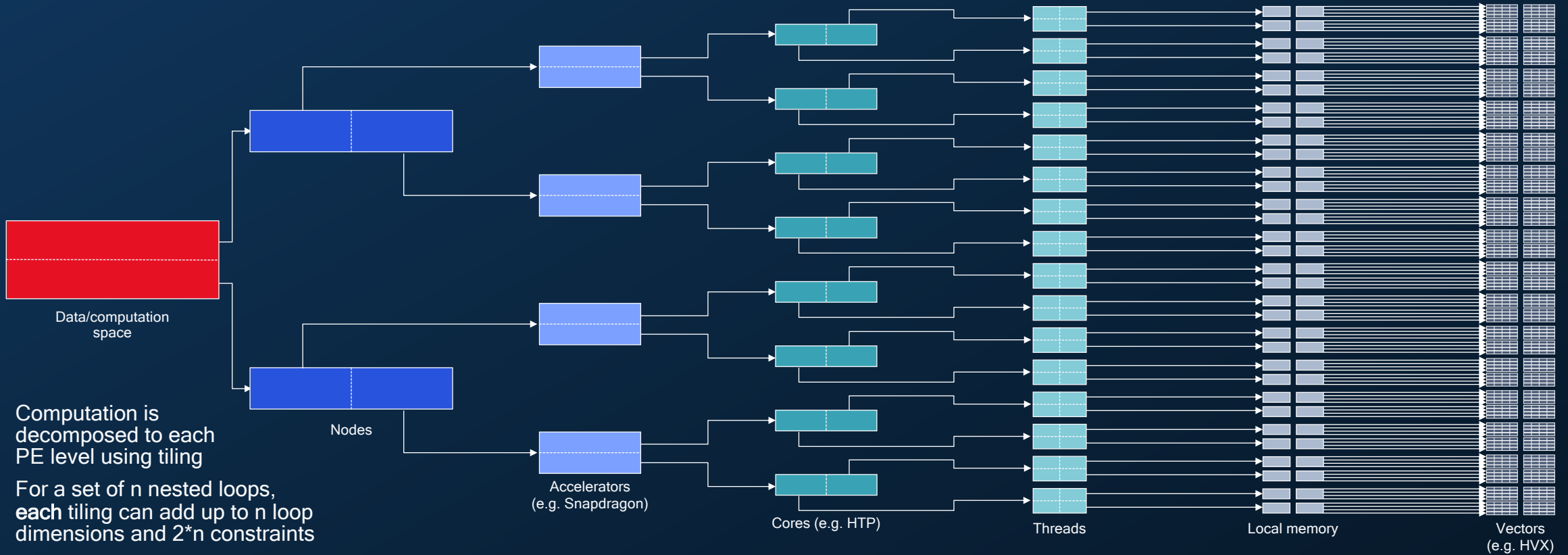
Hierarchical mapping optimization

Improving tractability of polyhedral optimization for deep PE hierarchies



Some evolutions of AI models are tractable in polyhedral AI compiler

We address the difficult ones and leverage the easy ones



for i_0 in $[1, n_{ti0}]$
 for j_0 in $[1, n_{tj0}]$

for i_1 in $[1, n_{ti1}]$
 for j_1 in $[1, n_{tj1}]$

for i_2 in $[1, n_{ti2}]$
 for j_2 in $[1, n_{tj2}]$

for i_3 in $[1, n_{ti3}]$
 for j_3 in $[1, n_{tj3}]$

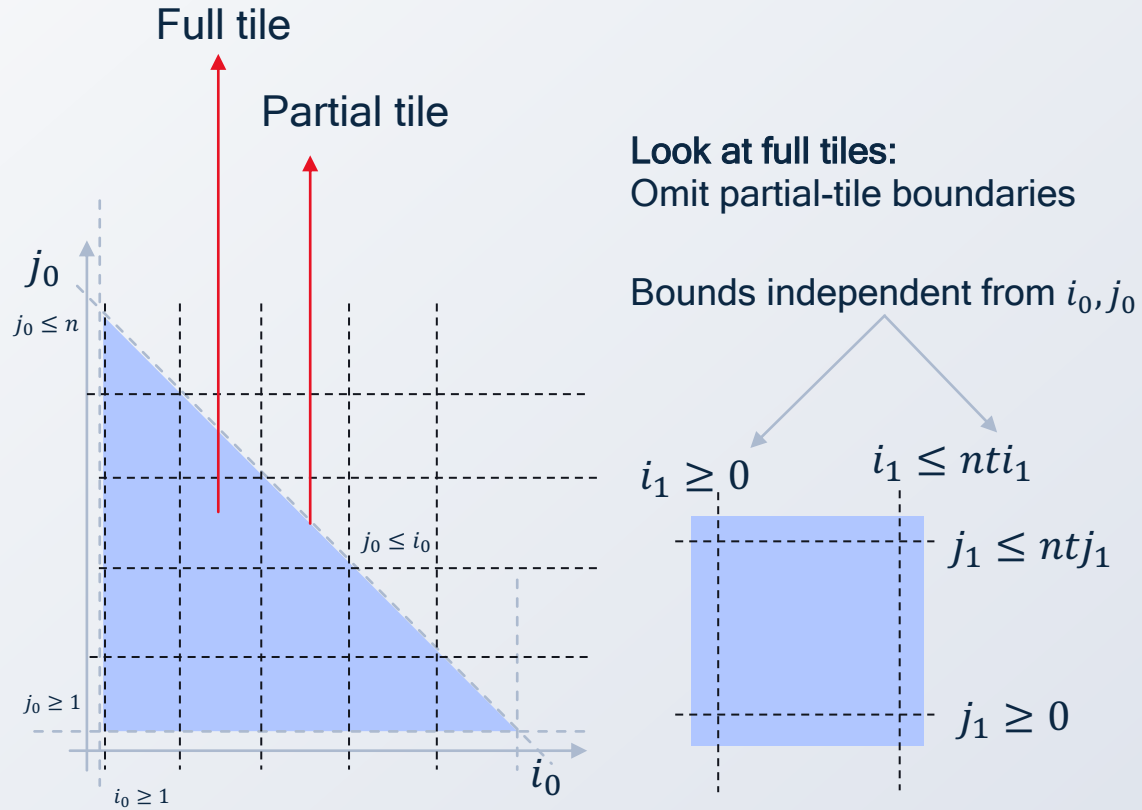
for i_4 in $[1, n_{ti4}]$
 for j_4 in $[1, n_{tj4}]$

for i_5 in $[1, n_{ti5}]$
 for j_5 in $[1, n_{tj5}]$

$f(i_0, j_0, \dots, i_5, j_5)$

Hierarchical mapping optimization
 addresses tractability issues to
 scale polyhedral compilation

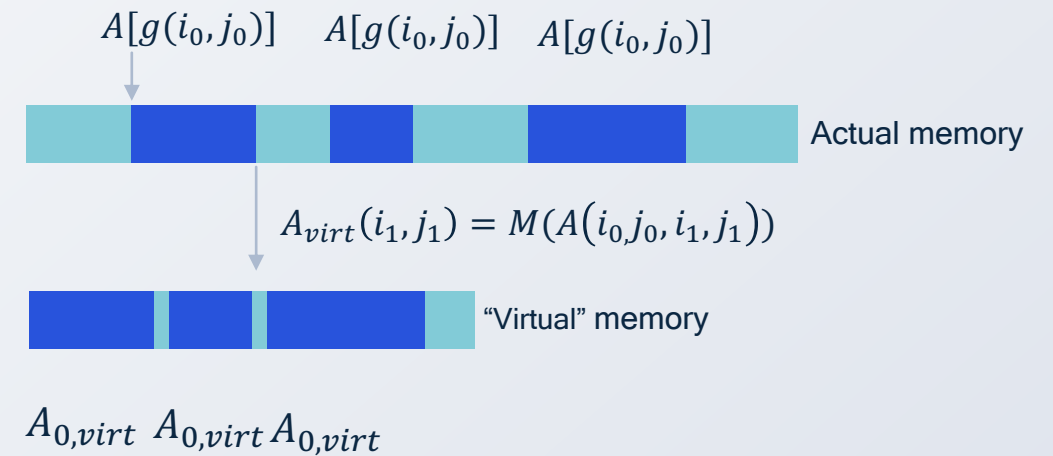
Loop transformations



[Pra16] Scalable hierarchical Polyhedral Compilation

Data transformations

Map tile data to a virtual space where accessed data regions are independent from other loop indices

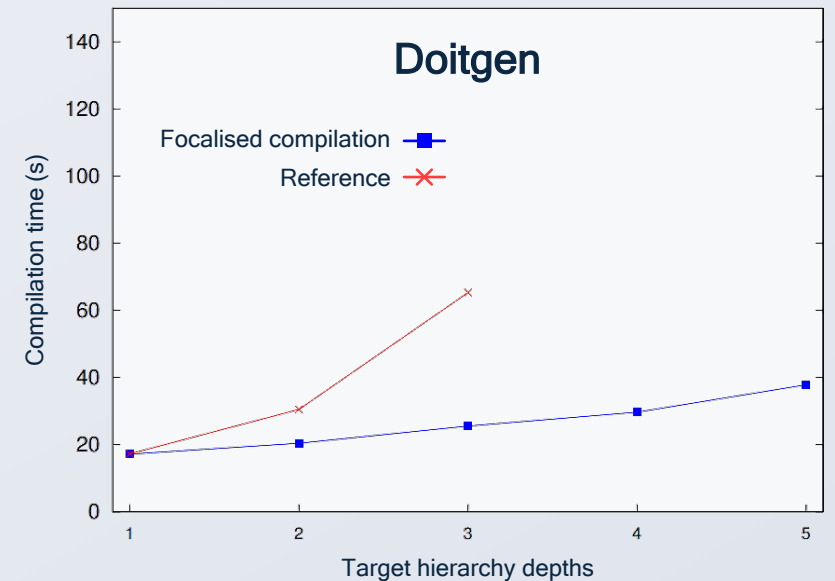
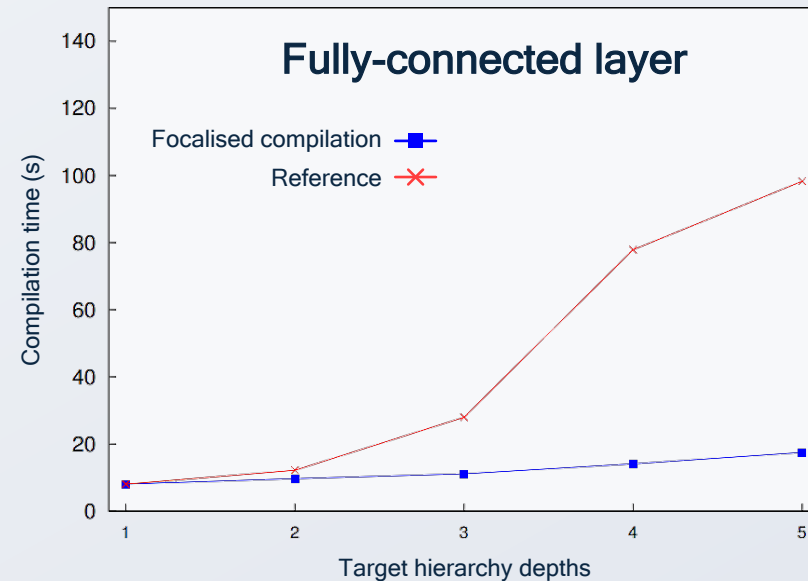
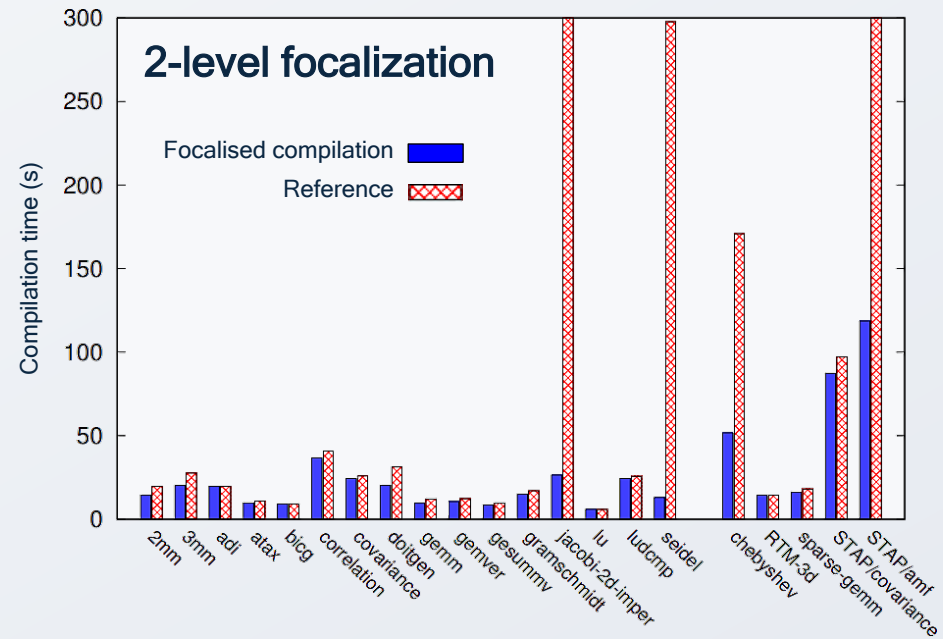


Focalization on a particular tiling level
with loop and data transformations

This makes intractable AI model evolutions tractable in AI compiler

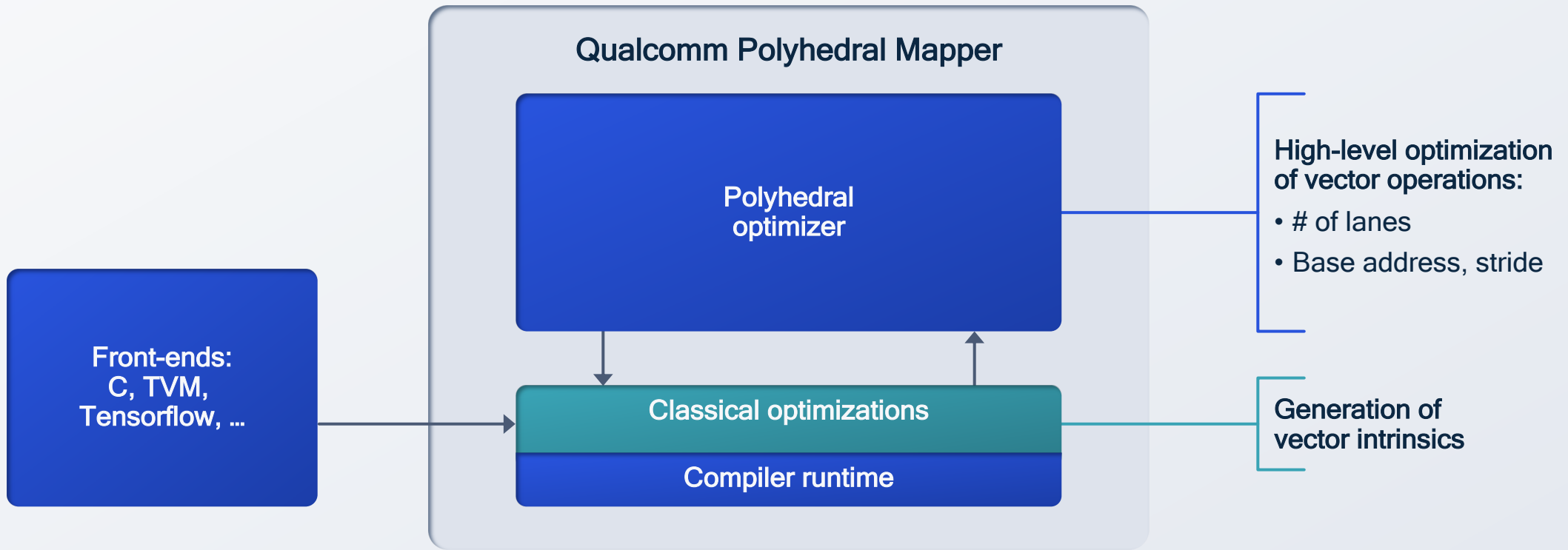
Using focalization, AI compile time reduces significantly

Successful on a broad class of programs



Auto-vectorization

In the Qualcomm Polyhedral Mapper



Separation of concerns:

Polyhedral Mapper finds independent operations suitable for vectorization

Underlying compiler generates optimized calls to vector intrinsics

Cooperative optimization between polyhedral and classical optimizations of the Qualcomm Polyhedral Mapper

Auto-optimized Transposed LayerNorm with the Qualcomm Polyhedral Mapper

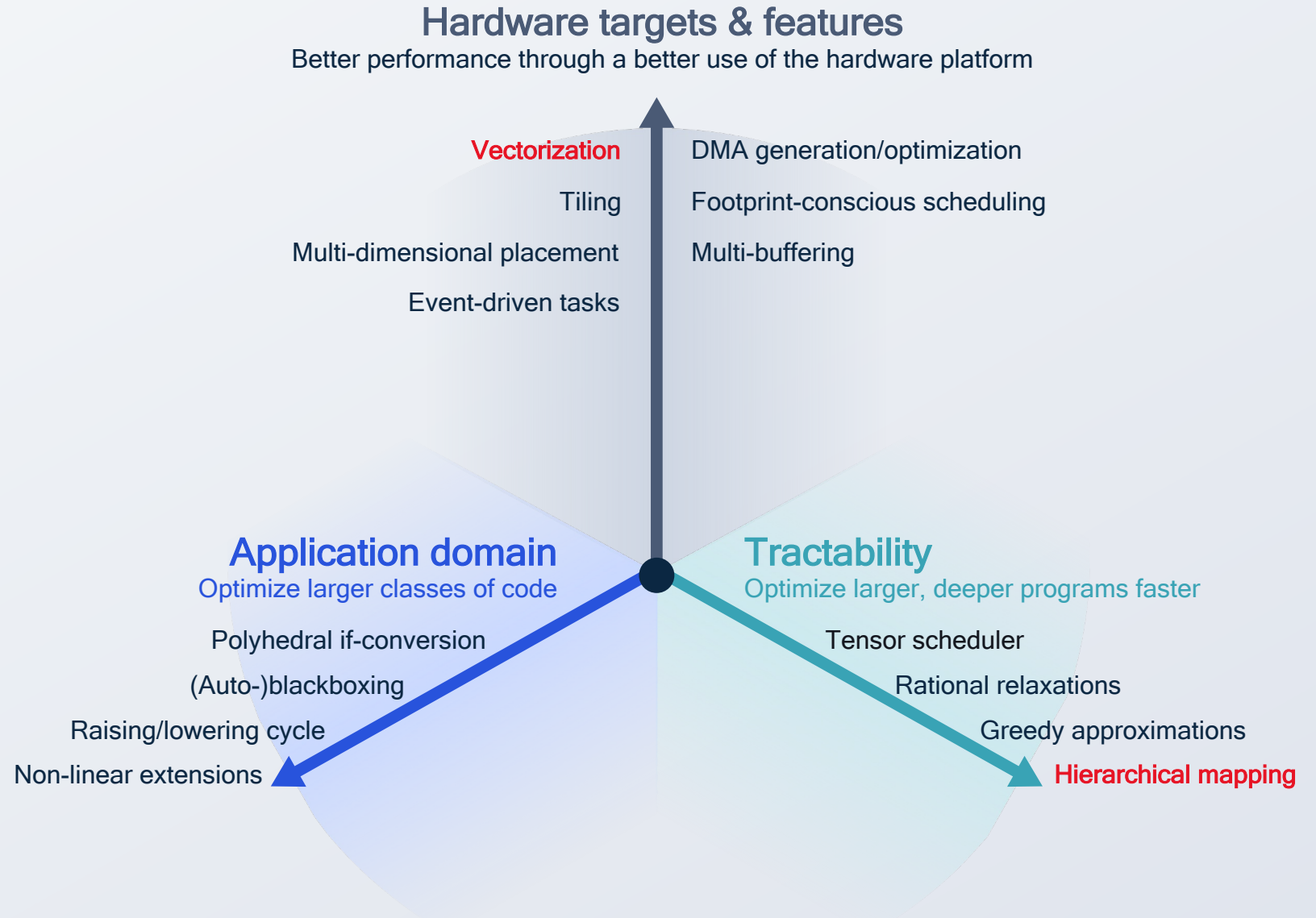
By automating
vectorization,
AI compiler allows
programmers to
write fewer lines
of code and
get substantial
performance

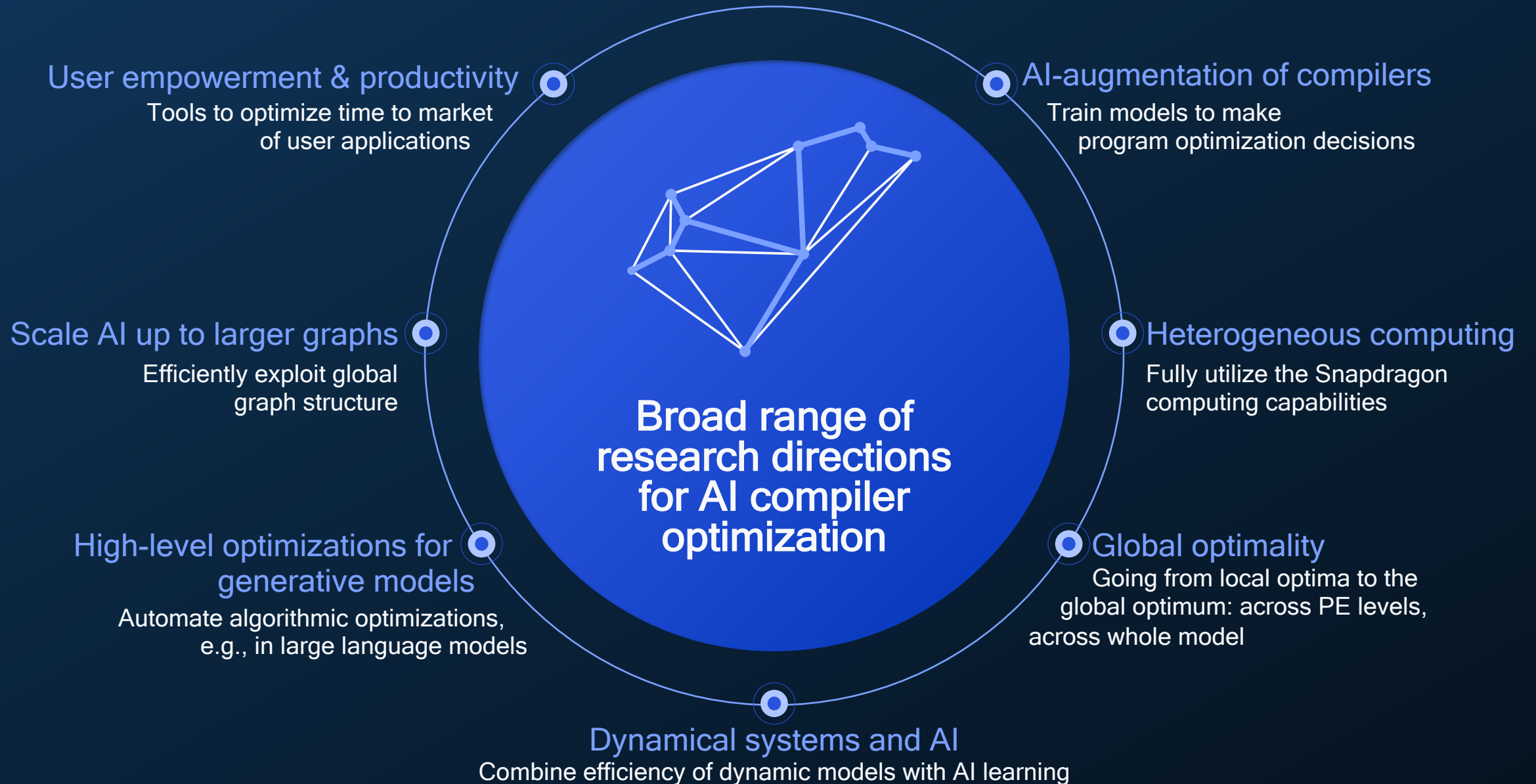
32.2% of hand-tuned
performance

620x speedup over
sequential performance

55.4 #optimized LOCs /
#original LOCs ratio

We are solving these challenges with polyhedral AI compilers

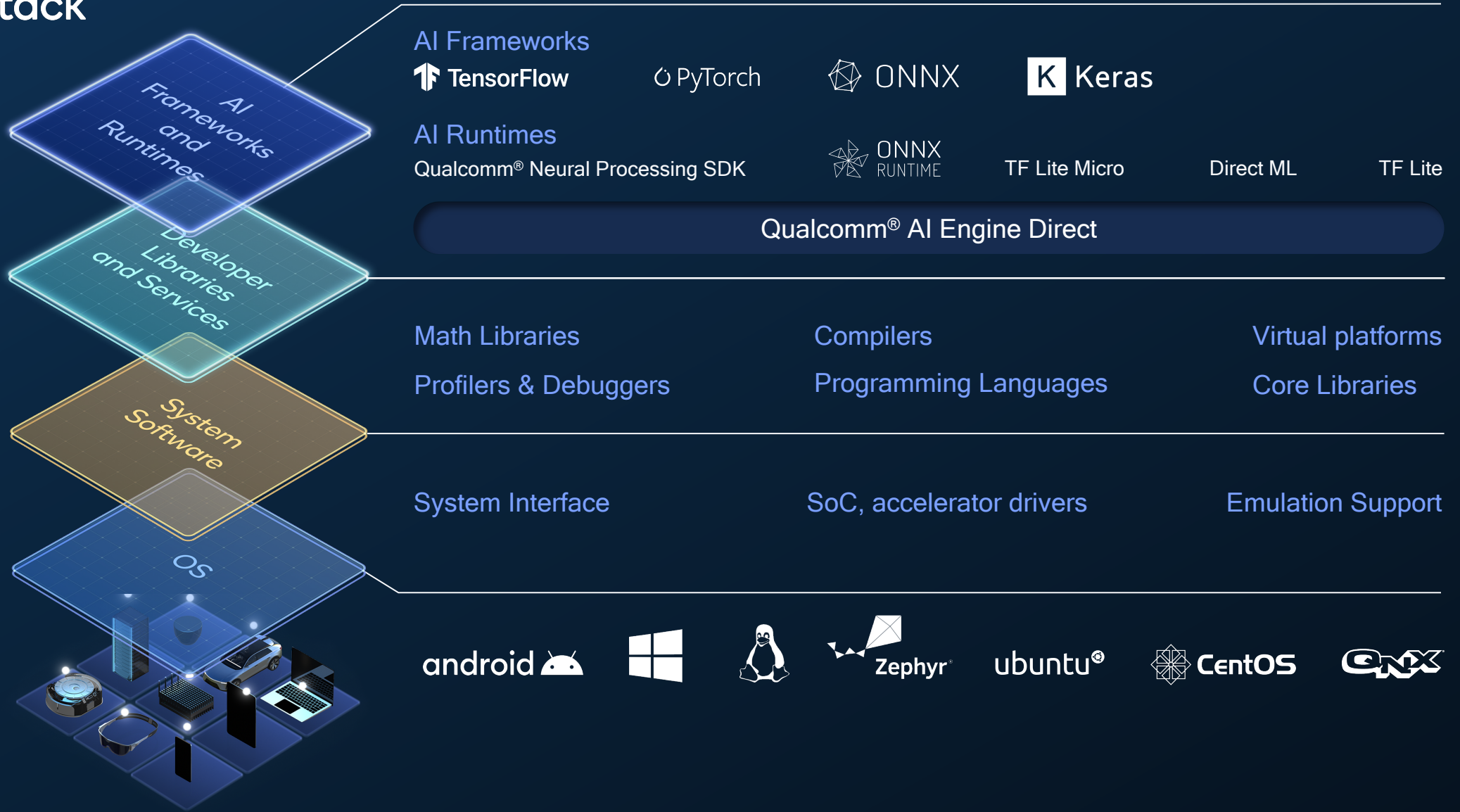






Qualcomm AI Stack

Qualcomm AI Studio



Qualcomm

Improved AI compilers are key to making machine learning ubiquitous

Qualcomm AI Research has achieved state-of-the-art results in AI compiler

We are enabling AI models to scale across a variety of use cases and run efficiently



Connect with us



www.qualcomm.com/research/artificial-intelligence



www.qualcomm.com/news/onq



[@QCOMResearch](https://twitter.com/QCOMResearch)



www.youtube.com/c/QualcommResearch



www.slideshare.net/qualcommwirelessevolution

Thank you

Qualcomm

Follow us on: [in](#) [twitter](#) [instagram](#) [youtube](#) [facebook](#)

For more information, visit us at:

qualcomm.com & qualcomm.com/blog

Nothing in these materials is an offer to sell any of the components or devices referenced herein.

©2018-2023 Qualcomm Technologies, Inc. and/or its affiliated companies. All Rights Reserved.

Qualcomm is a trademark or registered trademark of Qualcomm Incorporated. Other products and brand names may be trademarks or registered trademarks of their respective owners.

References in this presentation to “Qualcomm” may mean Qualcomm Incorporated, Qualcomm Technologies, Inc., and/or other subsidiaries or business units within the Qualcomm corporate structure, as applicable. Qualcomm Incorporated includes our licensing business, QTL, and the vast majority of our patent portfolio. Qualcomm Technologies, Inc., a subsidiary of Qualcomm Incorporated, operates, along with its subsidiaries, substantially all of our engineering, research and development functions, and substantially all of our products and services businesses, including our QCT semiconductor business.

Snapdragon and Qualcomm branded products are products of Qualcomm Technologies, Inc. and/or its subsidiaries. Qualcomm patented technologies are licensed by Qualcomm Incorporated.