

AISC: Approximate Instruction Set Computer

Alexandra Ferrerón¹, Jesús Alastruey-Benedé¹, Darío Suárez-Gracia¹, Ulya R. Karpuzcu²
¹Universidad de Zaragoza, Spain ²University of Minnesota, Twin Cities

Abstract

This paper makes the case for a single-ISA heterogeneous computing platform, AISC, where each compute engine (be it a core or an accelerator) supports a different subset of the very same ISA. An ISA subset may not be functionally complete, but the union of the (per compute engine) subsets renders a functionally complete platform-wide *single* ISA. Tailoring the microarchitecture of each compute engine to the subset of the ISA that it supports can easily reduce hardware complexity. At the same time, the energy efficiency of execution can improve by exploiting algorithmic noise tolerance: by mapping code sequences that can tolerate the incomplete ISA-subsets to the corresponding compute engines.

1 Motivation

The ISA specifies semantic and syntactic characteristics of a practically *functionally* complete set of machine instructions. Modern ISAs are not necessarily *mathematically* functionally complete, but provide sufficient expressiveness for practical algorithms. For software layers, the ISA defines the underlying machine – as capable as the variety of algorithmic tasks the composition of its building blocks, *instructions*, can express. For hardware layers, the ISA rather acts as a behavioral design specification for the machine organization. Accordingly, the ISA governs both the functional completeness and complexity of a machine design.

This paper makes the case for an alternative, single-ISA heterogeneous computing platform, AISC, which can reduce the ISA complexity, and thereby improve energy efficiency, on a per compute engine (be it a core or an accelerator) basis, without compromising the functional completeness of the overall platform. The distinctive feature of AISC is that each compute engine supports a *different subset* of the *very same* instruction set. Such per compute engine ISA subsets may be disjoint or overlapping. An ISA subset may not be functionally complete, but the union of the (per compute engine) subsets renders platform-wide a functionally complete *single* ISA. Therefore, software layers perceive AISC as a single-ISA machine. On the other hand, we can tailor the microarchitecture of each compute engine to the subset of the ISA that it supports. The result is less complex, more energy efficient compute engines, without compromising the overall functional completeness of the machine. To be able

to exploit this potential, we have to address many questions including

- Which subset of the ISA should each compute engine support?
- How to guarantee that each sequence of instructions scheduled to execute on a given compute engine only spans the respective subset of the ISA (with potential accuracy loss)? More specifically, how to map instruction sequences to the compute engines?
- How to keep the potentially incurred accuracy loss confined?

We can *approximate* the ISA per compute engine along two dimensions:

- *Horizontal* approximation simplifies instructions by reducing complexity (e.g., precision) on a per instruction basis. To be more specific, the subset of the ISA a compute engine implements in this case would selectively contain lower complexity (e.g., lower precision) instructions, by construction. Well-studied precision reduction approaches [2, 5, 6, 8, 10–12, 14, 16, 17] are directly applicable in this context. Reducing the operand width often enables simplification in the corresponding arithmetic operation, in addition to a more efficient utilization of the available communication bandwidth for data (i.e., operand) transfer.
- *Vertical* approximation eliminates complex and less frequently used instructions.

The combination of the two dimensions, *Vertical+Horizontal*, is also possible: In this case, the compute engine concerned would be able to *approximately* emulate complex and less frequently used instructions (that its ISA subset does not contain) by a sequence of simpler instructions. Along both dimensions, AISC trades computation accuracy for the complexity (and thereby, energy efficiency) on a per compute engine basis. The compiler and the runtime scheduler have to carefully choose compute engines in scheduling instructions to keep any potential accuracy loss below acceptable thresholds. At the same time, as the entire platform still supports the full-fledged ISA, instruction sequences not prone to approximation can still run at full accuracy.

AISC can also be regarded as an aggressive variant of architectural core salvaging [9] or ultra-reduced instruction set coprocessors [15], where actual hardware faults impair a compute engine’s capability to implement a subset of its ISA (and all compute engines support the same ISA by construction). Both of these studies detail how to achieve full-fledged functional completeness under the hardware-fault-induced loss of support for a subset of instructions. AISC, on the

other hand, features compute engines with approximate, i.e., incomplete, ISAs by construction.

2 Proof-of-concept Implementation

Let us start with a motivating example. Fig. 1 shows how the (graphic) output of a typical noise-tolerant application, SRR¹, changes for representative *Vertical*, *Horizontal*, and *Horizontal + Vertical* approximation under AISC. The application is compiled with GCC 4.8.4 with `-O1` on an Intel[®] Core™ i5 3210M machine. As we perform manual transformations on the code, high optimization levels hinder the task; we resort to `-O1` for our proof-of-concept and leave for future work more exploration on compiler optimizations. We focus on the main kernels where the actual computation takes place, and conservatively assume that this entire code would be mapped to a compute engine with approximated ISA. We use ACCURAX metrics [1] to quantify the accuracy-loss. We prototype basic *Horizontal* and *Vertical* ISA approximations on Pin 2.14 [7]. Fig. 1(a) captures the output for the baseline for comparison, *Native* execution, which excludes any approximation. We observe that the accuracy loss remains barely visible and varies under different approximations. Let us next take a closer look at the sources of this diversity.

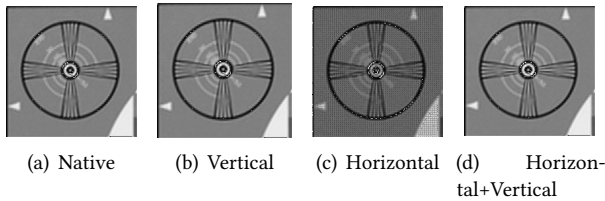


Figure 1. Graphic output of SRR benchmark under representative AISC approximations (b)-(d).

2.1 Vertical Approximation

The key question is how to pick the instructions to drop. A more general version of this question, *which instructions to approximate under AISC*, already applies across all dimensions, but the question becomes more critical in this case. As the most aggressive in our bag of tricks, *Vertical* can incur significant loss in accuracy. The targeted recognition-mining-synthesis applications can tolerate errors in data-centric phases as opposed to control [3]. Therefore, confining instruction dropping to data-flow can help limit the incurred accuracy loss. Fig. 1(b) captures an example execution outcome, where we randomly deleted static (arithmetic) floating point instructions. For each static instruction, we based the dropping decision on a pre-defined threshold t . We generated a random number r in the range $[0, 1]$, and dropped the static instruction if r remains below t . We experimented with threshold values between 1% and 10%.

¹Super Resolution Reconstruction, a computer vision application from the Cortex suite [13]. We use the (64×64) “EIA” input data set of 16 frames. The output is the (256×256) reconstructed image.

2.2 Horizontal Approximation

Without loss of generality, we experimented with three approximations to reduce operand widths: *DPtoSP*, *DP(SP)toHP*, and *DP(SP)toINT*. Under the IEEE 754 standard, 32 (64) bits express a single (double) precision floating point number: one bit specifies the *sign*; 8 (11) bits, the *exponent*; and 23 (52) bits the *mantissa*, i.e., the fraction. For example, $(-1)^{sign} \times 2^{exponent-127} \times 1.mantissa$ represents a single-precision floating number. *DPtoSP* is a *bit discarding* variant, which omits 32 least-significant bits of the mantissa of each double-precision operand of an instruction, and keeps the exponent intact. *DP(SP)toHP* comes in two flavors. *DPtoHP* omits 48 least-significant bits of the mantissa of each double-precision operand of an instruction, and keeps the exponent intact; *SPtoHP*, 16 least-significant bits of the mantissa of each single-precision operand of an instruction. Fig. 1(c) captures an example execution outcome under *DPtoHP*. *DP(SP)toINT* also comes in two flavors. *DPtoINT* (*SPtoINT*) replaces double (single) precision instructions with their integer counterparts, by rounding the floating point operand values to the closest integer.

2.3 Horizontal + Vertical Approximation

Without loss of generality, we experimented with two representatives in this case: *MULtoADD* and *DIVtoMUL*. *MULtoADD* converts multiplication instructions to a sequence of additions. We picked the smaller of the factors as the multiplier (which determines the number of additions), and rounded floating point multipliers to the closest integer number. *DIVtoMUL* converts division instructions to multiplications. We first calculated the reciprocal of the divisor, which next gets multiplied by the dividend to render the end result. In our proof-of-concept implementation based on the x86 ISA, the reciprocal instruction has 12-bit precision. *DIVtoMUL12* uses this instruction. *DIVtoMUL.NR*, on the other hand, relies on one iteration of the Newton-Raphson method [4] to increase the precision of the reciprocal to 23 bits. *DIVtoMUL12* can be regarded as an approximate version of *DIVtoMUL.NR*, eliminating the Newton-Raphson iteration, and hence enforcing a less accurate estimate of the reciprocal (of only 12 bit precision). Fig. 1(d) captures an example execution outcome under *DIVtoMUL.NR*.

References

- [1] I. Akturk, K. Khatamifard, and U. R. Karpuzcu. 2015. On Quantification of Accuracy Loss in Approximate Computing. In *12th Annual Workshop on Duplicating, Deconstructing and Debunking (WDDD)*.
- [2] V. K. Chippa, D. Mohapatra, K. Roy, S. T. Chakradhar, and A. Raghunathan. 2014. Scalable Effort Hardware Design. *IEEE Trans. on Very Large Scale Integration (VLSI) Systems* 22, 9 (Sept. 2014).
- [3] H. Cho, L. Leem, and S. Mitra. 2012. ERSA: Error Resilient System Architecture for Probabilistic Applications. *IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems* 31, 4 (April 2012).
- [4] M.D. Ercegovic and T. Lang. 2004. *Digital Arithmetic*. Morgan Kaufmann.

221	[5] Chih-Chieh Hsiao, Slo-Li Chu, and Chen-Yu Chen. 2013. Energy-aware Hybrid Precision Selection Framework for Mobile GPUs. <i>Computures and Graphics</i> 37, 5 (Aug. 2013).	276
222		277
223		278
224	[6] A. Jain, P. Hill, M.A. Laurenzano, M.E. Haque, M. Khan, S. Mahlke, L. Tang, and J. Mars. 2016. CPSA: Compute Precisely Store Approximately. In <i>Workshop on Approximate Computing Across the Stack</i> .	279
225		280
226	[7] Chi-Keung Luk, Robert Cohn, Robert Muth, Harish Patil, Artur Klauser, Geoff Lowney, Steven Wallace, Vijay Janapa Reddi, and Kim Hazelwood. 2005. Pin: Building Customized Program Analysis Tools with Dynamic Instrumentation. In <i>ACM SIGPLAN Conf. on Programming Language Design and Implementation</i> .	281
227		282
228		283
229		284
230	[8] T. Moreau, A. Sampson, L. Ceze, and M. Oskin. 2016. Approximating to the Last Bit. In <i>Workshop on Approximate Computing Across the Stack</i> .	285
231		286
232	[9] Michael D. Powell, Arijit Biswas, Shantanu Gupta, and Shubhendu S. Mukherjee. 2009. Architectural Core Salvaging in a Multi-core Processor for Hard-error Tolerance. In <i>ISCA</i> .	287
233		288
234		289
235	[10] Cindy Rubio-González, Cuong Nguyen, Hong Diep Nguyen, James Demmel, William Kahan, Koushik Sen, David H. Bailey, Costin Iancu, and David Hough. 2013. Precimonious: Tuning Assistant for Floating-point Precision. In <i>Int. Conf. on High Performance Computing, Networking, Storage and Analysis</i> .	290
236		291
237		292
238		293
239	[11] Adrian Sampson, Werner Dietl, Emily Fortuna, Danushen Gnanapragasam, Luis Ceze, and Dan Grossman. 2011. EnerJ: Approximate Data Types for Safe and General Low-power Computation. In <i>32nd ACM SIGPLAN Conf. on Programming Language Design and Implementation</i> .	294
240		295
241		296
242	[12] Mark Stephenson, Jonathan Babb, and Saman Amarasinghe. 2000. Bidwidth Analysis with Application to Silicon Compilation. In <i>ACM SIGPLAN Conf. on Programming Language Design and Implementation</i> .	297
243		298
244	[13] S. Thomas, C. Gohkale, E. Tanuwidjaja, T. Chong, D. Lau, S. Garcia, and M. Bedford Taylor. 2014. CortexSuite: A synthetic brain benchmark suite. In <i>IEEE Int. Symp. on Workload Characterization</i> .	299
245		300
246		301
247	[14] Ying Fai Tong, Rob A. Rutenbar, and David F. Nagle. 1998. Minimizing Floating-point Power Dissipation via Bit-width Reduction. In <i>Power-Driven Microarchitecture Workshop</i> .	302
248		303
249		304
250	[15] D. Wang, A. Rajendiran, S. Ananthanarayanan, H. Patel, M. V. Tripunitara, and S. Garg. 2014. Reliable Computing with Ultra-Reduced Instruction Set Coprocessors. <i>IEEE Micro</i> 34, 6 (2014).	305
251		306
252	[16] Thomas Y. Yeh, Glenn Reinman, Sanjay J. Patel, and Petros Faloutsos. 2009. Fool Me Twice: Exploring and Exploiting Error Tolerance in Physics-based Animation. <i>ACM Trans. on Graphics</i> 29, 1 (Dec. 2009).	307
253		308
254	[17] Q. Zhang, T. Wang, Y. Tian, F. Yuan, and Q. Xu. 2015. ApproxANN: An Approximate Computing Framework for Artificial Neural Network. In <i>Design, Automation & Test in Europe Conf. Exhibition</i> .	309
255		310
256		311
257		312
258		313
259		314
260		315
261		316
262		317
263		318
264		319
265		320
266		321
267		322
268		323
269		324
270		325
271		326
272		327
273		328
274		329
275		330