



HAL
open science

EOLE: Paving the Way for an Effective Implementation of Value Prediction

Arthur Perais, André Seznec

► **To cite this version:**

Arthur Perais, André Seznec. EOLE: Paving the Way for an Effective Implementation of Value Prediction. [Research Report] RR-8402, 2013, pp.25. hal-00907973v1

HAL Id: hal-00907973

<https://inria.hal.science/hal-00907973v1>

Submitted on 22 Nov 2013 (v1), last revised 29 Jan 2014 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



EOLE: Paving the Way for an Effective Implementation of Value Prediction

Arthur Perais, André Sez nec

**RESEARCH
REPORT**

N° 8402

November 2013

Project-Teams ALF

ISRN INRIA/RR--8402--FR+ENG

ISSN 0249-6399



EOLE: Paving the Way for an Effective Implementation of Value Prediction

Arthur Perais, André Seznec

Project-Teams ALF

Research Report n° 8402 — November 2013 — 25 pages

Abstract: Even in the multicore era, there is a continuous demand to increase the performance of single-threaded applications. However, the conventional path of increasing both issue width and instruction window size inevitably leads to the power wall. Value prediction (VP) was proposed in the mid 90's as an alternative path to further enhance the performance of wide-issue superscalar processors. Still, it was considered up to recently that a performance-effective implementation of Value Prediction would add tremendous complexity and power consumption in almost every stage of the pipeline.

Nonetheless, recent work in the field of VP has shown that given an efficient confidence estimation mechanism, prediction validation could be removed from the out-of-order engine and delayed until commit time. As a result, recovering from mispredictions via selective replay can be avoided and a much simpler mechanism – pipeline squashing – can be used, while the out-of-order engine remains mostly unmodified. Nonetheless, VP and validation at commit time entail strong constraints on the Physical Register File. Write ports are needed to write predicted results and read ports are needed in order to validate them at commit time, potentially rendering the overall number of ports unbearable. Fortunately, VP also implies that many single-cycle ALU instructions have their operands predicted in the front-end and can be executed in-place and in-order. Similarly, the execution of single-cycle instructions whose result has been predicted can be delayed until just before commit since predictions are validated at commit time.

Consequently, a significant number of instructions – 10% to 60% in our experiments – can bypass the out-of-order engine, allowing the reduction of the issue width, which is a major contributor to both out-of-order engine complexity and register file port requirement. This reduction paves the way for a truly practical implementation of Value Prediction. Furthermore, since Value Prediction in itself usually increases performance, our resulting {Early | Out-of-Order | Late} Execution architecture (EOLE), is often more efficient than a baseline VP-augmented 6-issue superscalar while having a significantly narrower 4-issue out-of-order engine.

Key-words: Microarchitecture, Value Prediction, VTAGE, EOLE

**RESEARCH CENTRE
RENNES – BRETAGNE ATLANTIQUE**

Campus universitaire de Beaulieu
35042 Rennes Cedex

EOLE: Paving the Way for an Effective Implementation of Value Prediction

Résumé : Même à l'ère des multicœurs, il existe une demande continue pour l'augmentation de la performance sur les applications mono-threads. Cependant, la solution conventionnelle consistant à augmenter la largeur d'exécution ainsi que la taille de la fenêtre d'instructions se heurte inévitablement au mur de la consommation. La Prédiction de Valeurs (VP) a été proposée dans les années 90 comme une alternative permettant d'améliorer la performance des processeurs superscalaires. Cela étant, une implémentation intéressante du point de vue cout-efficacité était jusqu'ici considérée comme impossible à cause de la complexité ainsi que de la consommation induite.

Cependant, des travaux récents dans le domaine de la Prédiction de Valeurs ont montré qu'avec un mécanisme d'estimation de la confiance efficace, la validation d'une prédiction pouvait être repoussée au moment où l'instruction est retirée du pipeline. Conséquemment, récupérer d'une mauvaise prédiction via une ré-exécution sélective peut-être évité et un mécanisme bien plus simple – vidage du pipeline – peut-être utilisé. Toute la partie du processeur chargée d'exécuter les instructions dans le désordre n'est donc pas modifiée.

Néanmoins, VP et la validation au retraitement impliquent des contraintes fortes sur le fichier de registres. Des ports d'écriture sont requis pour écrire les prédictions et des ports de lecture sont requis pour valider les prédictions au retraitement. Heureusement, VP implique aussi que beaucoup d'instructions simples ont leurs opérandes disponibles tôt dans le pipeline et peuvent être exécutées dans l'ordre. De façon similaire, l'exécution des instructions simples ayant été prédites peut être reportée aux derniers étages du pipeline puisque les prédictions sont validées au retraitement.

Au final, une proportion significative des instructions – 10% to 60% dans notre étude – peuvent contourner le moteur d'exécution dans le désordre, ce qui permet de réduire la largeur d'exécution, qui contribue grandement à la complexité du processeur. Cette réduction ouvre la porte à une implémentation réaliste de la Prédiction de Valeurs. De plus, puisque la VP augmente la performance, notre architecture {Early | Out-of-Order | Late} Execution architecture (EOLE), est souvent plus performante qu'une architecture superscalaire implémentant la VP tout en ayant un moteur d'exécution dans le désordre bien moins complexe.

Mots-clés : Microarchitecture, Prédiction de Valeurs, EOLE

1 Introduction & Motivations

Even in the multicore era, the need for higher single thread performance is driving the definition of new high-performance cores. Although the usual superscalar design does not scale, increasing the ability of the processor to extract *Instruction Level Parallelism* (ILP) by increasing the window size as well as the issue width has generally been the favored way to enhance sequential performance. For instance, consider the recently introduced Intel Haswell micro-architecture that has 33% more issue capacity than Intel Nehalem¹. To accommodate this increase, both the *Reorder Buffer* (ROB) and Scheduler size were substantially increased². On top of this, modern schedulers must support complex mechanisms such as *speculative scheduling* to enable back-to-back execution and thus *selective replay* to efficiently recover from schedule mispredictions [17].

In addition, the issue width impacts other structures: The *Physical Register File* (PRF) must provision more read/write ports as the width grows, while the number of physical registers must also increase to accommodate the ROB size. Because of this, both latency and power consumption increase and using a monolithic register file rapidly becomes complexity-ineffective. Similarly, a wide-issue processor should provide enough functional units to limit resource contention. Yet, the complexity of the bypass network grows quadratically with the number of functional units and quickly becomes critical regarding cycle time [23]. In other words, the out-of-order engine impact on power consumption and cycle time is ever increasing [6].

In this paper, we propose a modified superscalar design, the *{Early | Out-of-Order | Late} Execution* microarchitecture: *EOLE*. It is built on-top of a Value Prediction (VP) pipeline. VP allows dependents to issue earlier than previously possible by using predicted operands, and thus uncovers more ILP. Yet, predictions must be verified to ensure correctness. Fortunately, Perais and Seznec observed that one can validate the predicted results outside the out-of-order engine at commit time, provided an enhanced confidence estimation mechanism [24].

With *EOLE*, we leverage this observation to further reduce the complexity of the out-of-order (OoO) execution engine and the number of ports required on the PRF when VP is implemented. We achieve this reduction without significantly reducing overall performance. Our contribution is therefore twofold: First, *EOLE* paves the way to truly practical implementations of VP; second, it reduces complexity in the most complicated and power-hungry part of a modern OoO core.

In particular, when using VP, a significant number of single-cycle instructions have their operands ready in the front-end thanks to the value predictor. As such, we introduce *Early Execution* to execute simple instructions in-order in parallel with Rename by using predicted and/or immediate operands. Early-executed instructions are **not** sent to the OoO engine. Moreover, delaying VP validation until commit time avoids having to use *selective replay* and enforces a complete pipeline squash on a value misprediction. This guarantees that the operands of committed early-executed instructions are the correct operands. Early Execution requires simple hardware and reduces pressure on the OoO instruction window.

Similarly, since predicted results can be validated outside the OoO engine at commit time [24], we can offload the execution of predicted instructions to some dedicated in-order *Late Execution* pre-commit stage, where no *Select & Wakeup* has to take place. This does not hurt performance since instructions dependent on predicted instructions will simply use the predicted results rather than wait in the Scheduler. Similarly, the resolution of high confidence branches can be offloaded to the Late Execution stage since they are very rarely mispredicted. Overall, a total of 10% to 60% of the retired instructions can be offloaded from the OoO core.

As a result, *EOLE* benefits from both aggressiveness of modern OoO designs and higher

¹State-of-the-art in 2009

²From respectively 128 and 36 entries to 192 and 60 entries.

energy-efficiency of more conservative in-order designs. We evaluate EOLE against a baseline OoO model featuring VP and show that it achieves similar levels of performance having only 66% of the baseline issue capacity and a significantly less complex physical register file. This is especially interesting since it provides architects extra design headroom in the OoO engine to implement new architectural features.

The remainder of this paper is organized as follows. Section 2 discusses related work and provides some background on Value Prediction. Section 3 details the EOLE microarchitecture, which implements both *Early* and *Late Execution* by leveraging *Value Prediction*. Section 4 describes our simulation framework while Section 5 presents experimental results. Section 6 focuses on the qualitative gains in complexity and power consumption permitted by EOLE. Finally, Section 7 provides concluding remarks and directions for future research.

2 Related Work

Complexity-Effective Architectures Many propositions aim at reducing complexity in modern superscalar designs. In particular, it has been shown that most of the complexity and power consumption reside in the OoO engine, including the PRF [37], Scheduler and bypass network [23]. As such, previous studies focused on either reducing the complexity of existing structures, or in proposing new pipeline organizations.

Farkas et al. propose the *Multicluster* architecture in which execution is distributed among several execution clusters, each of them having its own register file [8]. The Alpha 21264 [15] is an example of real-world clustered architecture and shares many traits with the *Multicluster* architecture.

Palacharla et al. introduce a *dependence-based* microarchitecture where the centralized instruction window is replaced by several parallel FIFOs [23]. This greatly reduces complexity since only the head of each FIFO has to be selected by the *Select* logic. They also study a *clustered dependence-based* architecture to reduce the amount of bypass and window logic by using clustering.

Tseng and Patt propose the *Braid* architecture [35], which is fairly similar to the *clustered dependence-based* architecture except that instruction steering is done at compile time.

Austin proposes *Dynamic Implementation Validation* (DIVA) to check instruction results just before commit time, allowing the core to be faulty [2]. An interesting observation is that the latency of the checker has very limited impact on performance. This hints that adding pipeline stages between *Writeback* and *Commit* does not actually impact performance much.

Fahs et al. study *Continuous Optimization* where common compile-time optimizations are applied dynamically in the Rename stage [7]. This allows to *early execute* some instructions in the front-end instead of the OoO core. Similarly, Petric et al. propose RENO which also dynamically applies optimizations at rename-time [26].

Instead of studying new organizations of the pipeline, Kim and Lipasti present the *Half Price Architecture* [16]. They argue that many instructions are single operands and that both operands of dual operands instructions rarely become ready at the same time. Thus, the load capacitance on the tag broadcast bus can be greatly reduced by *sequentially waking-up* operands. Similarly, Ernst and Austin propose *Tag Elimination* to limit the number of comparators used for *Wakeup* [6].

Regarding the register file, Kim and Lipasti observe that many issuing instructions do not need to read both their operands in the register file since one or both can be available on the bypass network [16]. Thus, provisioning two read ports per issue slot is generally over-provisioning. Reducing the number of ports drastically reduces the complexity of the register file as ports are

much more expensive than registers.

Lastly, Lukefahr et al. propose to implement two back-ends – in-order and OoO – in a single core [20] and to dispatch instructions to the most adapted one. In most cases, this saves power at a slight cost in performance. In a sense, *EOLE* has similarities with such a design since instructions can be executed in different locations. However, no decision has to be made as the location where an instruction will be executed depends only on its type and status (e.g. predicted or not).

Note that our proposal is orthogonal to all these contributions since it reduces the number of instructions that enters the OoO execution engine.

Value Prediction *EOLE* builds upon the broad spectrum of research on Value Prediction independently initiated by Lipasti et al. and Gabbay et al. [10, 18].

Sazeides et al. refine the taxonomy of Value Prediction by categorizing predictors [28]. Specifically, they define two classes of value predictors: *Computational* and *Context-based*. On the one hand, *Computational* predictors generate a prediction by applying a function to the value(s) produced by the previous instance(s) of the instruction. For example, the Stride predictor [21] and the *2-Delta* Stride predictor [5] use the addition of a constant (stride).

On the other hand, *Context-Based* predictors rely on patterns in the value history of a given static instruction to generate predictions, e.g. the Finite Context Method (FCM) predictors [28]. Most of the initial studies on Value Prediction either assumed that the recovery on a value misprediction is immediate and induces – almost – no penalty [18, 19, 21, 39], or simply focused on accuracy and coverage rather than actual speedup [12, 22, 27, 28, 34, 38]. The latter studies were essentially ignoring the performance loss associated with misprediction recovery, i.e. assuming a perfect zero-cycle selective replay. Such a perfect selective replay mechanism is known to be unrealistic [17].

In a recent study, Perais and Seznec show that all value predictors are amenable to very high accuracy at the cost of some coverage [24]. This allows to delay prediction validation until commit time, thus removing the burden of implementing a complex selective replay mechanism. As such, the OoO engine remains mostly untouched by Value Prediction. In the same paper, they introduce the VTAGE context-based predictor. As the ITTAGE indirect branch predictor [30], VTAGE uses global branch history to select predictions.

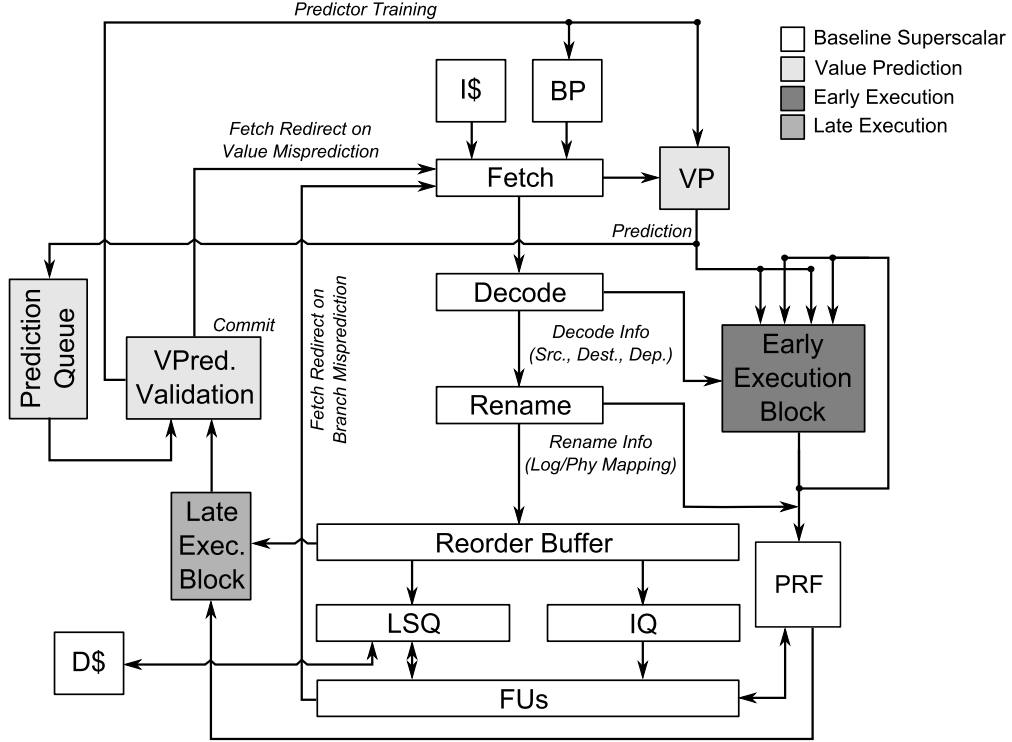
3 EOLE

3.1 Enabling EOLE Using Value Prediction

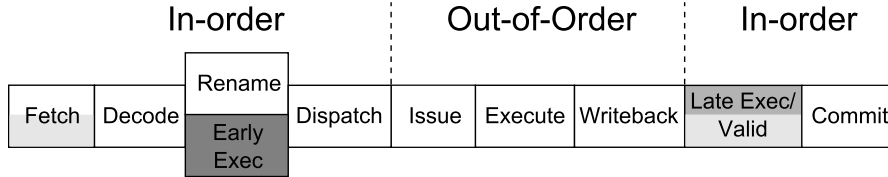
As described in the introduction, EOLE consists in adding a set of simple ALUs in the in-order front-end to early-execute instructions in parallel with Rename, and a second set in the in-order back-end to late-execute instructions just before they are committed.

While EOLE is heavily dependent on Value Prediction, they are in fact complementary features. Indeed, the former needs a value predictor to predict operands for Early Execution and provide temporal slack for Late Execution, while Value Prediction needs EOLE to reduce complexity and to become truly practical. However, EOLE requires prediction validation to be done at commit since 1) Validating at Execute forbids Late Execution and introduces strong constraints on the OoO engine (e.g. PRF port requirements) 2) Using *selective replay* to recover from a value misprediction nullifies the interest of both Early and Late Execution as all instructions must flow through the OoO window in case they need to be replayed.

To that extent, Perais and Seznec have proposed a confidence estimation mechanism greatly limiting the number of value mispredictions, *Forward Probabilistic Counters* (FPC) [24]. With



(a) Block diagram.



(b) Pipeline diagram.

Figure 1: The EOLE μ -architecture.

FPC, there is no need to validate as early as possible, and *pipeline squashing* at commit time can be used as the recovery mechanism. This enables the implementation of both Early and Late Execution, hence EOLE.

By eliminating the need to dispatch and execute many instructions in the OoO engine, EOLE substantially reduces the pressure on complex and power-hungry structures. As a result, those structures can be scaled down, yielding a less complex architecture whose performance is on-par with a more aggressive design. Moreover, doing so is orthogonal to previously proposed mechanisms such as *clustering* [8, 15, 23, 31] and does not *require* a centralized instruction window, even though this is the model we use to illustrate this paper. Fig. 1 depicts the EOLE architecture, implementing both Early Execution (darkest), Late Execution (darker) and Value Prediction (lighter). In the following paragraphs, we detail the two additional blocks required to implement EOLE and their interactions with the rest of the pipeline.

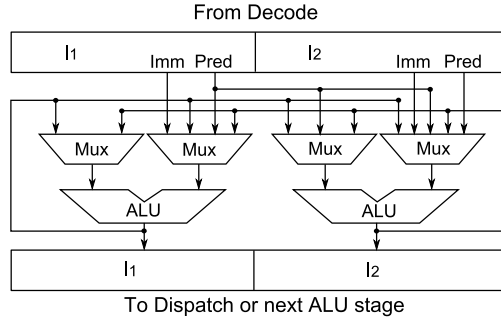


Figure 2: Early Execution Block. The logic controlling the ALUs and muxes is not shown for clarity.

3.2 Early Execution Hardware

The core idea of Early Execution is to position one or more ALU stages in the front-end in which instructions with available operands will be executed. For complexity concerns, however, it seems necessary to limit Early Execution to single-cycle ALU instructions. Indeed, implementing complex functional units in the front-end to execute multi-cycle instructions does not appear as a worthy tradeoff. Early Execution is done in-order, hence, it does not require renamed registers and can take place in parallel with Rename. For instance, Fig. 2 depicts the Early Execution Block adapted to a 2-wide Rename stage.

Renaming is often pipelined over several cycles. Consequently, we can use several ALU stages and simply insert pipeline registers between each stage. The actual execution of an instruction can then happen in any of the ALU stages, depending on the readiness of its operands coming from *Decode* (i.e. immediate), the local³ bypass network (i.e. from instructions early-executed in the previous cycle) or the value predictor. Operands are **never** read from the PRF.

In a nutshell, all eligible instructions simply flow through the ALU stages, propagating their results in each bypass network accordingly once they have executed. Finally, at the end of the last stage, results as well as predictions are written to the PRF.

An interesting design concern lies with the number of stages required to capture an interesting proportion of instructions. However, we actually found that using more than a single stage was highly inefficient, as illustrated in Fig. 3. This figure shows the proportion of committed instructions eligible for Early Execution for a baseline 8-wide rename, 6-issue model (see Table 1 in Section 4), using the VTAGE/2D-Stride hybrid predictor later described in Table 2, Section 4. As a result, in further experiments, we consider a 1-deep Early Execution Block.

To summarize, Early Execution only requires a single new computing block, which is shown in dark grey in Fig. 1. In particular, the mechanism we propose does not require any storage area to store temporaries as all values are living inside the pipeline registers or the bypass network. Finally, since we execute in-order, each instruction is mapped to a single ALU and scheduling is straightforward.

³For complexity concerns, we consider that bypass does not span several stages. Consequently, if an instruction depends on a result computed by an instruction located two rename-groups ahead, it will not be early-executed.

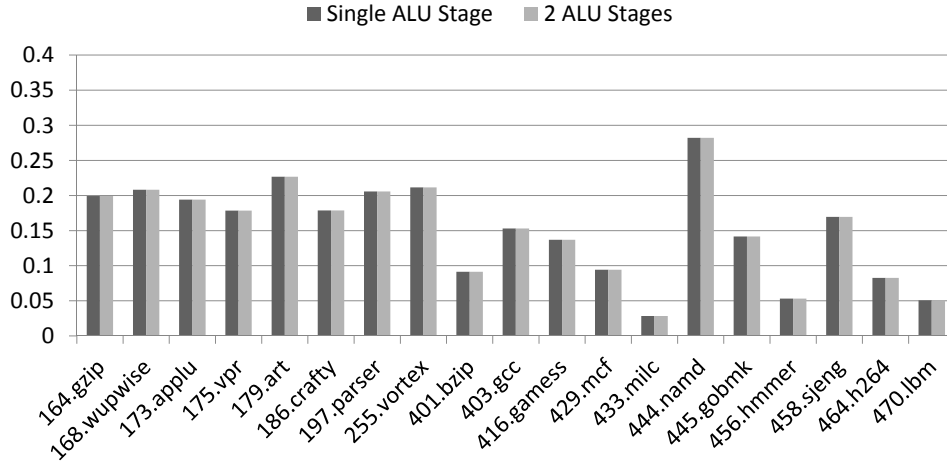


Figure 3: Proportion of committed instructions that can be early-executed, using one or two ALU stages and a VTAGE-2DStride hybrid predictor (later described in Section 4).

3.3 Late Execution Hardware

Late Execution targets instructions whose results have been predicted⁴. It is done just before validation time, that is, out of the execution engine. As for Early Execution, we limit ourselves to single-cycle ALU instructions to minimize complexity. Interestingly, Seznec showed in [29] that conditional branch predictions flowing from TAGE can be categorized such that very high confidence predictions are known. Since these high confidence branches exhibit a misprediction rate generally lower than 0.5%, resolving them in the Late Execution block has a marginal impact on overall performance. Consequently, we consider both single-cycle predicted ALU instructions and very high confidence branches⁵ for Late Execution. Furthermore, predicted instructions can also be early-executed. In that event, they only need to be validated in case another early-executed instruction from the same rename-group used the prediction as an operand.

In this study, we did not try to set confidence on the other branches (indirect jumps, returns). Yet, provided a similar high confidence estimator for these categories of branches, one could postpone the resolution of high confidence ones until the Late Execution stage.

In any case, Late Execution further reduces pressure on the OoO engine in terms of instructions dispatched to the Scheduler. As such, it also removes the need for predicting only critical instructions [9, 27, 36] since minimizing the number of instructions flowing through the OoO engine requires maximizing the number of predicted instructions. Hence, usually useless predictions from a performance standpoint become useful in EOLE. Fig. 4 shows the proportion of committed instructions that can be late-executed using a baseline 6-issue processor with a VTAGE-2DStride hybrid predictor (respectively described in Tables 1 and 2 in Section 4).

Late Execution needs to implement *commit width* ALUs and the associated read ports in the PRF. In particular, even if an instruction I_1 to be late-executed depends on the result of instruction I_0 also to be late-executed (both are in the same commit-group), it does not need to wait as it can use the predicted result of I_0 . In other words, all non executed instructions reaching the Late Execution stage have all their operands ready, as in DIVA [2]. Due to the need to validate predictions as well as late-execute some instructions, at least one extra pipeline stage

⁴Instructions eligible for prediction are μ -ops producing a result that can be read by a subsequent μ -op, as defined by the ISA implementation.

⁵Predictions whose *pred* counter is saturated [29].

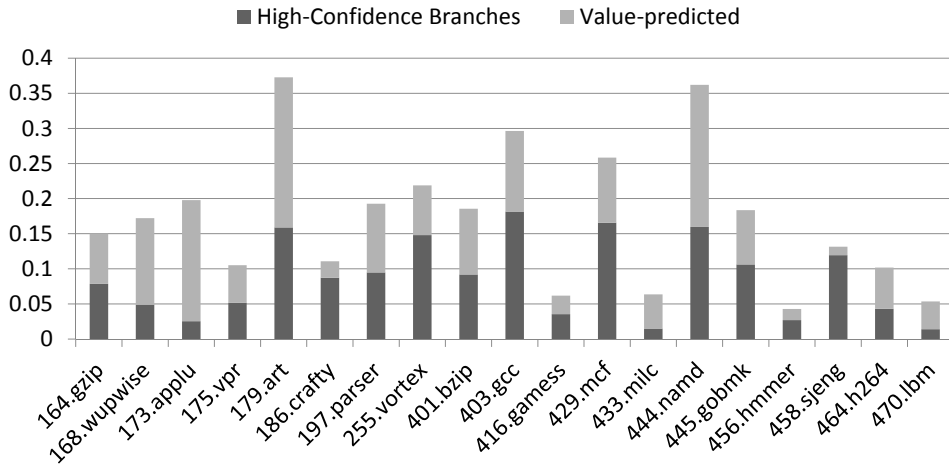


Figure 4: Proportion of committed instructions that can be late-executed using a VTAGE-2DStride (see Section 4) hybrid predictor. Late-executable instructions that can also be early-executed are not counted since instructions are executed once at most.

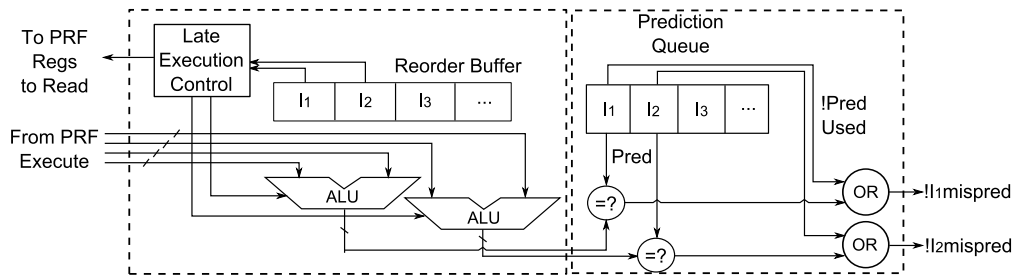


Figure 5: Late Execution Block for a 2-wide processor. The left part can late-execute two instructions while the right part validates two results against their respective predictions. Buses are *register-width*-bit wide.

after *Writeback* is likely to be required in EOLE.

Consequently, the hardware needed for Late Execution is simple. It does not even require a bypass network. A high-level view of a 2-wide Late Execution Block is shown in Fig. 5. In further experiments, we consider that Late Execution and prediction validation can be done in the same cycle, before the *Commit* stage. EOLE is therefore only one cycle longer than the baseline superscalar it is compared to. While this may be optimistic due to the need to read from the PRF, this only impacts the value misprediction penalty and the pipeline fill delay. In particular, since low confidence branches are resolved in the same cycle as for the baseline, the average branch misprediction penalty will remain very similar. Lastly, as a first step, we also consider that enough ALUs are implemented (i.e. as many as the commit-width). As a second step, we shall consider reduced-width Late Execution.

3.4 Potential OoO Engine Offload

We obtain the ratio of retired instructions that can be offloaded from the OoO engine for each benchmark by summing the columns in Fig. 3 and 4. This ratio is very dependent on the

application, ranging from less than 10% for *milc*, *hmmmer* and *lbm* to more than 50% for *art* and up to 60% for *namd*. Nonetheless, it represents a significant part of the retired instructions in most cases.

4 Evaluation Methodology

4.1 Simulator

In our experiments, we use the x86_64 ISA to validate EOLE, even though EOLE can be adapted to any general-purpose ISA. We use a modified⁶ version of the *gem5* cycle-accurate simulator [3]. However, contrarily to modern x86 implementations, *gem5* does not support *move elimination* [7, 14, 26], *μ -op fusion* [11] and does not implement a *stack-engine* [11]. We consider a relatively aggressive 4GHz, 6-wide issue superscalar⁷ baseline with a fetch-to-commit latency of 19 cycles. Since we focus on the OoO engine complexity, both in-order front-end and in-order back-end are overdimensioned to treat up to 8 μ -ops per cycle. We model a deep front-end (15 cycles) coupled to a shallow back-end (3 cycles) to obtain a realistic branch/value misprediction penalty. Table 1 describes the characteristics of the baseline pipeline we use in more details. In particular, the OoO scheduler is dimensioned with a unified centralized 64-entry IQ and a 192-entry ROB on par with Haswell’s, the latest commercially available Intel microarchitecture. We refer to this baseline as the *Baseline_6_64* configuration.

As μ -ops are known at Fetch in *gem5*, all the widths given in Table 1 are in μ -ops, even for the fetch stage. Independent memory instructions (as predicted by the Store Sets predictor [4]) are allowed to issue out-of-order. Entries in the IQ are released upon issue since *selective replay* is not needed on value mispredictions.

In the case where Value Prediction is used, we add a pre-commit stage responsible for validation (and Late Execution in EOLE). This accounts for an additional pipeline cycle (20 cycles) and an increased value misprediction penalty (21 cycles min.). Minimum branch misprediction latency remains unchanged except for mispredicted very high confidence branches when EOLE is used.

4.2 Value Predictor Operation

The predictor makes a prediction at Fetch for every eligible μ -op (i.e. producing a register that can be read by a subsequent μ -op, as defined by the ISA implementation). To index the predictor, we XOR the PC of the x86_64 instruction left-shifted by two with the μ -op number inside the x86_64 instruction. This avoids all μ -ops mapping to the same entry. We assume that the predictor can deliver as many predictions as requested by the Fetch stage.

In previous work, a prediction is written into the PRF and replaced by its non-speculative counterpart when it is computed in the OoO engine [24]. In parallel, predictions are put in a FIFO queue to be able to validate them – in-order – at commit time. In EOLE, we also use a queue for validation. However, instead of directly writing predictions to the PRF, we place predictions in the Early Execution units, which will in turn write the predictions to the PRF. By doing so, we can use predictions as operands in the Early Execution units.

⁶Our modifications mostly lie with the ISA implementation. In particular, we implemented branches with a single μ -op instead of three and we removed some false dependencies existing between instructions due to the way flags are renamed/written.

⁷On our benchmark set and with our baseline simulator, an 8-issue machine achieves only marginal speedup over this baseline.

Front End	L1I 4-way 32KB, Perfect TLB; 8-wide fetch (2 taken branch/cycle), decode, rename; TAGE 1+12 components [30] 15K-entry total, 20 cycles min. mis. penalty; 2-way 4K-entry BTB, 32-entry RAS;
Execution	192-entry ROB, 64-entry IQ unified, 48/48-entry LQ/SQ, 256/256 INT/FP registers; 1K-SSID/LFST Store Sets [4]; 6-issue, 6ALU(1c), 4MulDiv(3c/25c*), 6FP(3c), 4FPMulDiv(5c/10c*), 4Ld/Str; Full bypass; 8-wide retire;
Caches	L1D 4-way 32KB, 2 cycles, 64 MSHRs, 4 load ports; Unified L2 16-way 2MB, 12 cycles, 64 MSHRs, no port constraints, Stride prefetcher, degree 8, distance 1; All caches have 64B lines and LRU replacement;
Memory	Single channel DDR3-1600 (11-11-11), 2 ranks, 8 banks/rank, 8K row-buffer, tREFI 7.8us; Across a 64B bus; Min. Read Lat.: 75 cycles, Max. 185 cycles.

Table 1: Simulator configuration overview. *not pipelined.

Predictor	#Entries	Tag	Size (KB)
2D-Stride [5]	8192	Full (51)	251.9
VTAGE [24]	8192 (Base)	-	68.6
	6×1024	$12 + rank$	64.1

Table 2: Layout Summary. For VTAGE, *rank* is the position of the tagged component and varies from 1 to 6, 1 being the component using the shortest history length.

x86 Flags In the x86_64 ISA, some instructions write flags based on their results while some need them to execute (e.g. branches) [13]. We assume that flags are computed as the last step of Value Prediction, based on the predicted value. In particular, the *Zero Flag* (ZF), *Sign Flag* (SF) and *Parity Flag* (PF) can be easily inferred from the predicted result. Remaining flags – *Carry Flag* (CF), *Adjust Flag* (AF) and *Overflow Flag* (OF) – depend on the operands and cannot be inferred from the predicted result only. We found that always setting the *Overflow Flag* to 0 did not cause many mispredictions and that setting CF if SF was set was a reasonable approximation. The *Adjust Flag*, however, cannot be set to 0 or 1 in the general case. This is a major impediment to the value predictor coverage since we consider a prediction as incorrect if one of the derived flags – thus the flag register – is wrong. Fortunately, x86_64 forbids the use of decimal arithmetic instructions. As such, AF is not used and we can simply ignore its correctness when checking for a misprediction [13].

Predictors Considered in this Study Previous work in the field of Value Prediction has provided us with many predictors. In this study, we do not explore all possibilities. Rather, we focus on the hybrid predictor VTAGE-2DStride predictor introduced by Perais and Sez nec [24]. It combines a simple and cost-effective *2-Delta* Stride predictor [5] as a representative of the *computational* family – as defined by Sazeides et al. [28] – and a state-of-the-art VTAGE predictor [24] as a representative of the *context-based* family. For confidence estimation, we use Forward Probabilistic Counters as described by Perais and Sez nec in [24]. In particular, we use 3-bit confidence counters whose forward transitions are controlled by the vector $v = \{1, \frac{1}{32}, \frac{1}{32}, \frac{1}{32}, \frac{1}{32}, \frac{1}{64}, \frac{1}{64}\}$ as we found it to perform best with VTAGE-2DStride.

Table 2 summarizes the configuration of each predictor component. If the resulting sizes are large, one has to consider that there exist many degrees of freedom impacting the effective size of the predictor components, such as the number of entries, the tag-length, the stride-width or

the number of components in VTAGE. However, finding the most cost-efficient configuration is beyond the scope of this paper, thus, we use large tables to demonstrate the potential of EOLE.

4.3 Benchmarks

Program	Input	IPC
164.gzip (INT)	input.source 60	0.984
168.wupwise (FP)	wupwise.in	1.553
173.applu (FP)	applu.in	1.591
175.vpr (INT)	net.in arch.in place.out dum.out -nodisp -place_only -init_t 5 -exit_t 0.005 - alpha_t 0.9412 -inner_num 2	1.326
179.art (FP)	-scanfile c756hel.in -trainfile1 a10.img - trainfile2 hc.img -stride 2 -startx 110 - starty 200 -endx 160 -endy 240 -objects 10	1.211
186.crafty (INT)	crafty.in	1.769
197.parser (INT)	ref.in 2.1.dict -batch	0.544
255.vortex (INT)	lendian1.raw	1.781
401.bzip2 (INT)	input.source 280	0.888
403.gcc (INT)	166.i	1.055
416.gamess (FP)	cytosine.2.config	1.929
429.mcf (INT)	inp.in	0.105
433.milc (FP)	su3imp.in	0.459
444.namd (FP)	namd.input	1.860
445.gobmk (INT)	13x13.tst	0.766
456.hmmmer (INT)	nph3.hmm	2.477
458.sjeng (INT)	ref.txt	1.321
464.h264ref (INT)	foreman_ref_encoder_baseline.cfg	1.312
470.lbm (FP)	reference.dat	0.748

Table 3: Benchmarks used for evaluation. Top: CPU2000, Bottom: CPU2006. INT: 12, FP: 7, Total: 19.

We use a subset of the the SPEC’00 [32] and SPEC’06 [33] suites to evaluate our contributions as we focus on single-thread performance. Specifically, we use 12 integer benchmarks and 7 floating-point programs⁸. Table 3 summarizes the benchmarks we use as well as their input, which are part of the *reference* inputs provided in the SPEC software packages. To get relevant numbers, we identify a region of interest in the benchmark using Simpoint 3.2 [25]. We simulate the resulting slice in two steps: First, warm up all structures (caches, branch predictor and value predictor) for 50M instructions, then collect statistics for 100M instructions.

5 Experimental Results

In our experiments, we first use *Baseline_6_64* as the baseline to gauge the impact of adding a value predictor only. Then, in all subsequent experiments, we use said baseline augmented with the predictor presented in Table 2 (*Baseline_VP_6_64*). Our objective is to characterize the potential of EOLE at decreasing the complexity of the OoO engine. We assume that Early and Late Execution stages are able to treat any group of up to 8 consecutive μ -ops every cycle. In Section 6, we will consider tradeoffs to enable realistic implementations.

⁸We do not use the whole suites due to some currently missing system calls in *gem5-x86*.

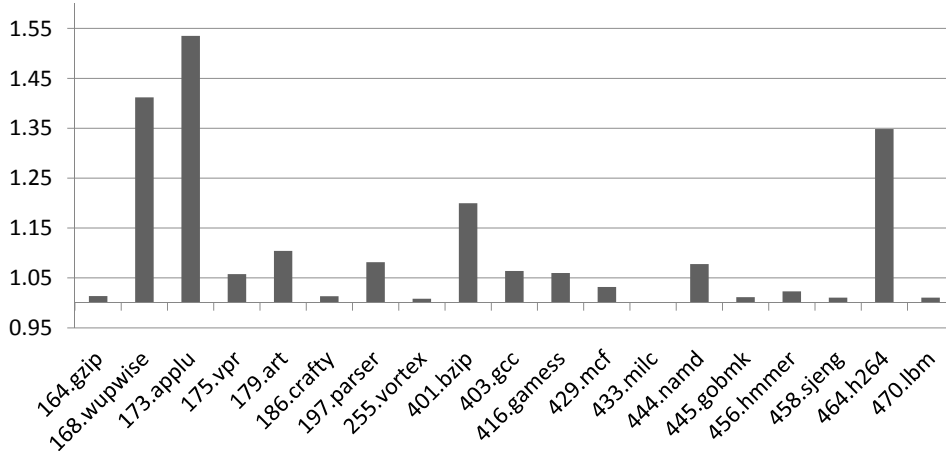


Figure 6: Speedup over *Baseline_6_64* brought by Value Prediction using VTAGE-2DStride as the predictor.

5.1 Performance of Value Prediction

Fig. 6 illustrates the performance benefit of augmenting the baseline processor with the VTAGE-2DStride predictor. A few benchmarks present interesting potential e.g. *wupwise*, *applu*, *bzip*, *h264*, some a more moderate potential e.g. *vpr*, *art*, *gamess*, *gcc*, *namd* and a few others low potential. No slowdown is observed.

In further experiments, illustration in performance figures will be speedups over the baseline described in Table 1, featuring the VTAGE-2DStride value predictor of Table 2. We refer to it as *Baseline_VP_6_64* configuration.

5.2 Issue Width Impact on EOLE

Applying EOLE without modifying the OoO core (*EOLE_6_64*) is illustrated in Fig. 7. In this figure, we also illustrate the baseline, but with a 4-way only issue OoO engine (*Baseline_VP_4_64*) and EOLE using a 4-way issue OoO engine (*EOLE_4_64*).

By itself, EOLE slightly increases performance over the baseline, with a few benchmarks achieving 5% speedup or higher. The particular case of *namd* is worth to be noted as with VP, it would have benefited from an 8-issue core by more than 10%. Through EOLE, we actually increase the number of instructions that can be executed at the same time, hence performance goes up in this benchmark.

Shrinking the issue width to 4 reduces the performance of the baseline by a significant factor on many applications, e.g. *applu*, *crafty*, *vortex*, *namd*, *hmmmer*, *sjeng* for which slowdown is more than 5% (up to 14% for *namd*). For EOLE, such a shrink only reduces performance by a few percent compared with *EOLE_6_64*. Furthermore, *EOLE_4_64* still performs slightly higher than *Baseline_VP_6_64* in several benchmarks e.g. *applu*, *vortex* and *namd*. A single slowdown of 1.8% is reported for *hmmmer*.

Therefore, EOLE can be considered as a mean to reduce issue width without impacting performance on a processor featuring VP.

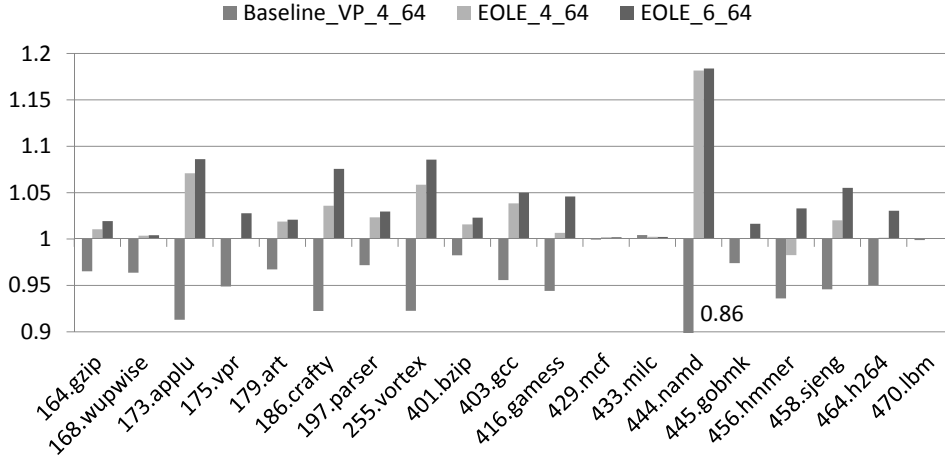


Figure 7: Performance of EOLE and the baseline with regard to issue width, normalized to *Baseline_VP_6_64*.

5.3 Impact of Instruction queue size

In Fig. 8, we illustrate the respective performance of shrinking the instruction queue size from 64 to 48 entries for the baseline and EOLE.

A few benchmarks are quite sensitive to such a shrink for *Baseline_VP_6_48* e.g. *wupwise*, *applu*, *crafty*, *vortex*, *namd*, *hmmr* and *h264*. On the other hand, *EOLE_VP_6_48* does not always exhibit the same behavior. Most applications encounter only minor losses with *EOLE_6_48* (less than 5% except for *hmmr*) and higher losses with *Baseline_6_48*, e.g. *applu* with 4% speedup against 9% slowdown, or *namd* with 16% speedup against 10% slowdown.

In practice, the benefit of EOLE is greatly influenced by the proportion of instructions that are not sent to the OoO engine. For instance *namd* needs a 64-entry IQ in the baseline case, but since it is also an application for which many instructions are predicted or early-executed, it can deal with a smaller IQ in the EOLE case.

On the other hand, *hmmr*, the application that suffers the most from reducing the instruction queue size with EOLE, exhibits a relatively low coverage of predicted or early-executed instructions. Nonetheless, with *EOLE_6_48*, slowdown is limited to 5% at most for all but one benchmark, *hmmr*, for which slowdown is around 9%.

5.4 Summary

EOLE provides opportunities for either slightly improving the performance over a VP-augmented processor without increasing the complexity of the OoO engine, or reaching the same level of performance with a significantly reduced OoO engine complexity. In the latter case, reducing the issue width is our favored direction as it addresses scheduler complexity, PRF complexity and bypass complexity. EOLE also mitigates the performance loss associated with a reduction of the instruction queue size.

In the next section, we provide directions to limit the global hardware complexity and power consumption induced by the EOLE design and the overall integration of VP in a superscalar processor.

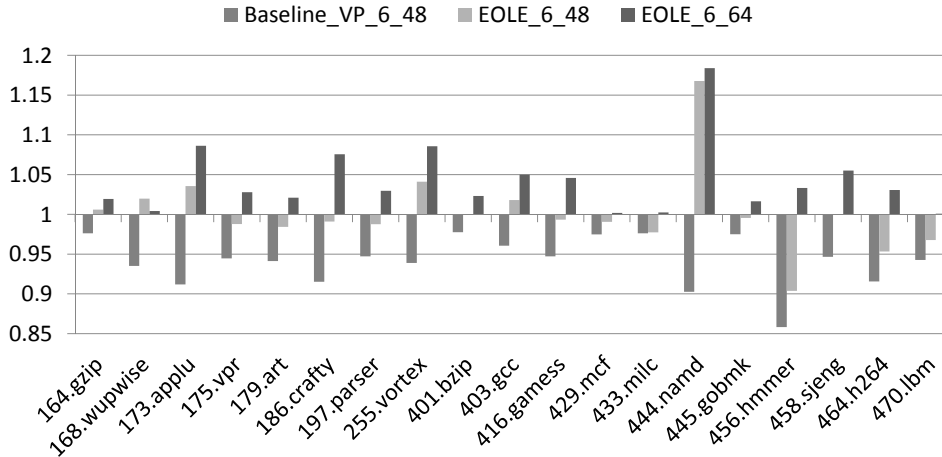


Figure 8: Performance of EOLE and the baseline with regard to the number of entries in the IQ, normalized to *Baseline_VP_6_64*.

6 Hardware Complexity

In the previous section, we have shown that, provided that the processor already implements Value Prediction, adopting the EOLE design may allow to use a reduced-issue OoO engine without impairing performance. On the other hand, extra complexity and power consumption are added in the Early Execution engine as well as the Late Execution engine.

In this section, we first describe the potential hardware simplifications on the OoO engine enabled by EOLE. Then, we describe the extra hardware cost associated with the Early Execution and Late Execution engines. Finally, we provide directions to mitigate this extra cost. Note however that a precise evaluation would require a complete processor design and is beyond the scope of this paper.

6.1 Shrinking the Out-of-Order Engine

Out-of-Order Scheduler Our experiments have shown that with EOLE, the OoO issue width can be reduced from 6 to 4 without significant performance loss on our benchmark set. This would greatly impact *Wakeup* since the complexity of each IQ entry would be lower. Similarly, a narrower issue width mechanically simplifies *Select*. As such, both steps of the *Wakeup & Select* critical loop could be made faster and/or less power hungry.

Providing a way to reduce issue width with no impact on performance is also crucial because modern schedulers must support complex features such as *speculative scheduling* and thus *selective replay* to recover from scheduling mispredictions [17].

Lastly, to our knowledge, most scheduler optimizations proposed in the literature can be added on top of EOLE. This includes the *Sequential Wakeup* of Kim and Lipasti [16] or the *Tag Elimination* of Ernst and Austin [6]. As a result, power consumption and cycle time could be further decreased.

Functional Units & Bypass Network As the number of cycles required to read a register from the PRF increases, the bypass network becomes more crucial, as it allows instructions to "catch" their operands as they are produced and thus execute back-to-back. However, a full

bypass network is very expensive, especially as the issue width – hence the number of functional units – increases. Ahuja et al. showed that partial bypassing could greatly impede performance, even for a simple in-order single-issue pipeline [1]. Consequently, in the context of a wide-issue OoO superscalar with a multi-cycle register read, missing bypass paths may cripple performance even more.

EOLE allows to reduce the issue width in the OoO engine. Therefore, it reduces the design complexity of a full bypass by reducing the number of ALUs and the number of simultaneous writers on the network.

A Limited Number of Register File Ports on the OoO Engine Through reducing the issue width on the OoO engine, EOLE mechanically reduces the number of read and write ports required on the PRF for regular OoO execution.

6.2 Extra Hardware Complexity Associated with Late/Early Execution

Cost of the Late Execution Block The extra hardware complexity associated with Late Execution consists of three main components. First, for validation at commit time, a prediction queue (FIFO) is required to store predicted results. This component is needed anyway as soon as VP associated with validation at commit time is implemented. Second, ALUs are needed for late execution. Last, the operands for the late-executed instructions or the computed results for predicted but not late-executed instructions must be read from the PRF, and results must be checked against the predicted results.

In the simulations presented in Section 5, we have assumed that up to 8 μ -ops (i.e. *commit-width*) could be late-executed per cycle. This would necessitate 8 ALUs and up to 16 read ports on the PRF.

Cost of the Early Execution Block A single stage of simple ALUs is sufficient to capture most of the potential benefit of Early Execution. The main hardware cost associated with Early Execution is this stage of ALUs and the associated full bypass. Additionally, the predicted results must be written on the register file.

Therefore, in our case, a complete 8-wide Early Execution stage necessitates 8 ALUs, a full 8-to-8 bypass network and 8 write ports on the PRF.

The Physical Register File From the above analysis, an EOLE-enhanced core featuring a 4-issue OoO engine (*EOLE_4_64*) would have to implement a PRF with a total of 12 write ports (resp. 8 for Early Execution and 4 for OoO execution) and 24 read ports (resp. 8 for OoO execution and 16 for Late Execution).

The area cost of a register file is approximately proportional to $(R + W) * (R + 2W)$, R and W respectively being the number of read and write ports [40]. That is, at equal number of registers, the area cost of the EOLE PRF would be 4 times the initial area cost of the 6-issue baseline (*Baseline_6_64*) PRF. Moreover, this would also translate in largely increased power consumption and longer access time, thus impairing cycle time and/or lengthening the register file access pipeline.

Without any optimization, *Baseline_VP_6_64* would necessitate 14 write ports (resp. 8 to write predictions and 6 for the OoO engine) and 20 read ports (resp. 8 for validation and 12 for the OoO engine), i.e. slightly less than *EOLE_4_64*. In both cases, this overhead might be considered as prohibitive in terms of silicon area, power consumption and access time.

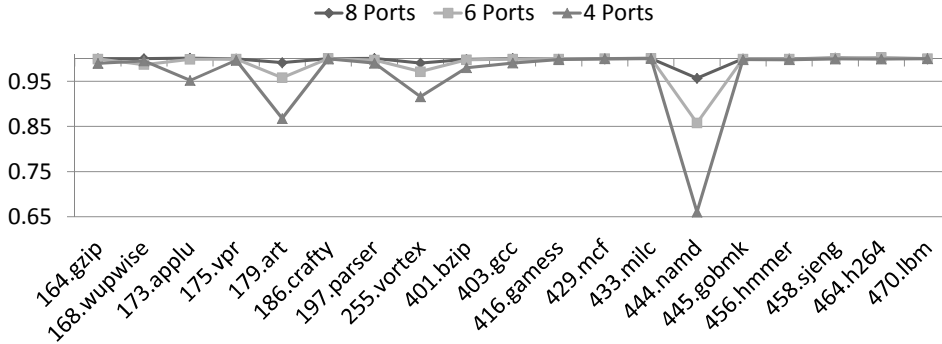


Figure 9: Performance of *EOLE_4_64* when the number of read ports dedicated to Late Execution/validation is limited, normalized to *EOLE_4_64* with 16 ports dedicated to LE/validation.

However, simple solutions can be devised to reduce the overall cost of the PRF and the global hardware cost of Early/Late Execution without significantly impacting the global performance. These solutions apply for EOLE as well as for a baseline implementation of VP. We describe said solutions below.

6.3 Mitigating the Hardware Cost of Early/Late Execution

Narrow Late Execution Not all instructions are predicted or late-executable (i.e. predicted and simple ALU or high confidence branches). Moreover, entire groups of 8 μ -ops are rarely ready to commit. Therefore, one can limit the number of potentially late-executed instructions and/or predicted instructions per cycle. For instance, the maximum commit-width can be kept to 8 with the extra constraint of using only 6 or 8 PRF read ports for Late Execution and validation.

In Fig. 9, we illustrate simulations assuming *EOLE_4_64* with respectively 4, 6 and 8 PRF read ports dedicated to prediction validation and late execution. As expected, the performance loss compared with unconstrained Late Execution is marginal with 8 PRF read ports, except for *namd* (around 4.5%). This is due to *namd* having a very large proportion of late-executed instructions, explaining the need for additional ports for Late Execution and validation. Carefully throttling the predictor is a possibility to reduce the impact of limited ports in benchmarks where the coverage of the predictor is very high (around 80% of the eligible instructions in *namd*). To summarize, 8 PRF ports for prediction validation and late execution are sufficient and except for *namd*, 6 can be envisioned.

Therefore, one can implement Late Execution and value prediction validation at commit on our baseline core with only 8 extra read ports on the PRF instead of 16. This does not significantly decrease performance.

Mitigating Early-Execution Hardware Cost There are several opportunities to limit the effective number of write ports on the PRF required by Early-Execution and Value Prediction.

A first observation is that many μ -ops are not predicted. Hence, they do not generate any writes on the PRF and as for Late Execution, one could limit the number of μ -ops that write to the PRF at the exit of the Early Execution stage. This was also suggested for Value Prediction only in [24]. In particular, after the Early Execution/Rename stage, the μ -ops and their predicted/computed results could be buffered. Dispatch groups of up to 8 instructions

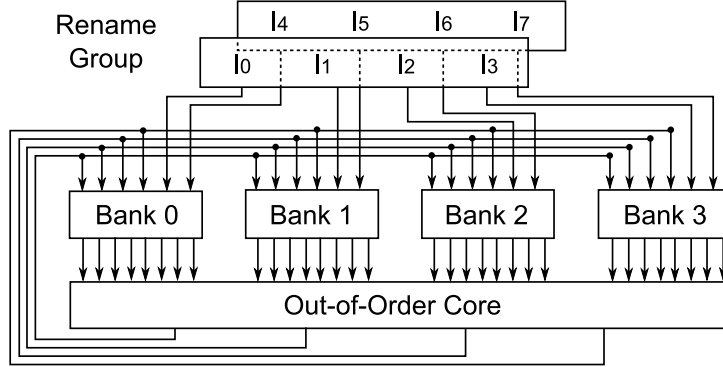


Figure 10: Organization of a 4-bank PRF supporting 8-wide Early Execution and 4-wide OoO issue. The additional read ports per-bank required for Late Execution/validation are not shown for clarity.

would be built, with the extra constraint of a limited number of at most 4 Early Execute writes on the PRF per dispatch group.

A second observation is that because instructions in the front-end are treated in-order and are therefore consecutive, one can use a banked PRF and force the allocation of physical registers for the same dispatch group to different register banks. For instance, considering a 4-bank PRF, out of a group of 8 consecutive μ -ops, 2 could be allocated to each bank. A dispatch group of 8 consecutive μ -ops would at most write 2 registers in a single bank after Early Execution. Thus, Early Execution would necessitate only two extra write ports on each PRF bank, as illustrated in Fig. 10 for an 8-wide Rename/Early Execute, 4-issue OoO core.

The two propositions above could be mixed to further limit the number of extra write port to only one. Logic to determine which μ -instruction will be predicted/early-executed should provide the necessary input to the rename logic.

Another possibility would be to specialize physical registers at allocation with the Early Executed or predicted μ -ops writing in a set of registers and the OoO engine writing in another distinct set of registers.

As the exploration of all these solutions would require a complete paper, we only illustrate the banked PRF and leave other mechanisms for future work. In particular, registers from distinct banks are allocated to consecutive μ -ops and Rename is stalled if the current bank does not have any free register. We consider respectively 2 banks of 128 registers, 4 banks of 64 registers and 8 banks of 32 registers⁹. As illustrated in Fig. 11, the performance loss associated with load unbalancing is quite limited for our benchmark set, and using 4 banks of 64 registers instead of a single bank of 256 registers appears as a reasonable tradeoff.

Lastly, note that in any case, very similar optimizations are needed for a practical implementation of a core featuring Value Prediction without EOLE.

The Overall Complexity of the Register File The elements presented above show that one could implement an EOLE core featuring a 4-way OoO engine with a PRF of reasonable complexity.

Assuming a reduced-issue (e.g. 4) EOLE core, 16 read ports (resp. 8 for the OoO engine and 8 for limited Late Execution/validation) would be needed. Depending on the mitigating

⁹8 banks were simulated. However, there could be situations where the whole set of architectural registers would be allocated to a single bank, leading to major functionality issues.

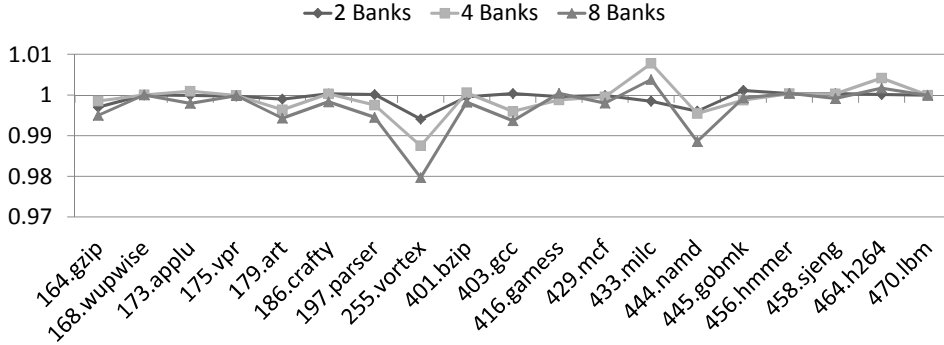


Figure 11: Performance of *EOLE_4_64* using a different number of banks in the PRF, normalized to *EOLE_4_64* with a single bank.

technique(s) chosen for Early Execution, the number of effective PRF write ports may vary. Limiting register writes originating from Early Execution would result in 8 write ports: 4 for Early Execution and 4 for the OoO engine. Banked register allocation would result in a 4-bank PRF featuring only 6 write ports per bank, 2 for Early Execution and 4 for OoO execution, as illustrated in Fig. 10. Combining limited register writes from Early Execution and banked register allocation could further reduce the number of write ports per PRF bank to 5.

Moreover, it should be mentioned that the EOLE structure should naturally lead to a distributed register file structure with one register file servicing reads from the OoO engine and the other register file servicing reads from the Late Execution engine. The PRF could be naturally built with two copies of a 6 write/8 read ports, 4-bank register file (or four copies of a 6 write/4 read ports, 4-bank register file).

On the other hand, if only Value Prediction (but validation at commit time) is implemented, the PRF is much more complex since a 6-issue OoO core is needed to ensure good performance. Register file ports are needed for read at validation time and writes at prediction time, but the PRF optimizations described for EOLE can be used. Assuming a 4-bank PRF – as for EOLE –, this would lead to each bank featuring 8 write ports (2 for predicted μ -ops and 6 for the OoO engine) and 16 read ports (4 read ports for validation assuming validating at most 4 predictions and 12 read ports for the out-of-execution order engine). Distributing the PRF would lead to a copy with 8 write ports and 12 read ports for the OoO engine (or 3 copies of a 8 write/4 read ports register file) and a copy with 8 write/4 read ports register file for the prediction validation stage. Fig. 12 illustrates performance for configurations with a limited number of ports and a 4-bank PRF.

Assuming, that the PRF is designed using either 6 write/4 read ports or 8 write/4 read ports building blocks, the previously mentioned area cost formula [40] allows us to roughly compare the respective relative PRF silicon areas in the three considered designs. The total area cost of the PRF for our target EOLE is $\frac{4}{3}$ that of the area cost of the baseline 6-way (no VP) PRF, but only $\frac{2}{3}$ that of the VP-enhanced 6-way baseline with banked allocation and reduced validation width. Moreover, as previously mentioned, one can use a distributed register file structure with EOLE. In that case, the PRF copy in the OoO engine would be only $\frac{2}{3}$ the area of the PRF of the baseline *without* Value Prediction. That is, the register file in the OoO engine is less likely to become a temperature hotspot than in a conventional design.

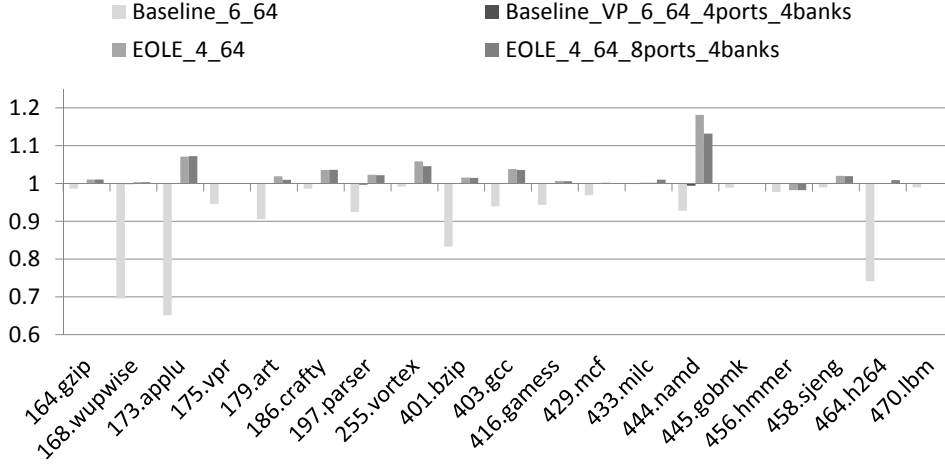


Figure 12: Performance of *EOLE_4_64* using 8 ports for Late Execution and validation and having 4 64-register banks, *EOLE_4_64* with 16 ports for LE/validation and a single bank and *Baseline_6_64*, normalized to *Baseline_VP_6_64*.

6.4 Summary

Apart from the prediction tables and the update logic, the major hardware overhead associated with implementing VP and VP validation at commit time comes from the extra read and write ports on the register file [24]. We have shown above that EOLE allows to get rid of most of this overhead on the PRF.

In particular, EOLE allows to use a 4-issue OoO engine instead of an 6-issue engine. This implies a much smaller instruction scheduler, a much simpler bypass network and a reduced number of PRF read and write ports in the OoO engine. As a result, one can expect many advantages in the design of the OoO execution core: Significant silicon area savings, significant power savings in the scheduler and the register file and savings on the access time of the register file. Power consumption savings are crucial since the scheduler has been shown to consume almost 20% of the power of a modern superscalar core [6], and is often a temperature hotspot in modern designs. As such, even if global power savings were not to be achieved due to the extra hardware required in EOLE, the power consumption will be more distributed across the core.

On the other hand, EOLE requires some extra but relatively simple hardware for Early/Late Execution. Apart from some relatively simple control logic, this extra hardware consists of a set of ALUs and a bypass network in the Early Execution stage and a set of ALUs in the Late Execution stage. A full rank of ALUs is actually unlikely to be needed. From Fig. 3, we presume that a rank of 4 ALUs would be sufficient.

Furthermore, implementing EOLE will not impair cycle time. Indeed, Early Execution requires only one stage of simple ALUs and can be done in parallel with Rename. Late Execution and validation may require more than one additional pipeline stage compared to a conventional superscalar processor, but this should have a fairly small impact since low-confidence branch resolution is not delayed. In fact, since EOLE simplifies the OoO engine, it is possible that the core could actually be clocked higher, yielding even more sequential performance.

Therefore, our claim is that EOLE makes a clear case for implementing VP on wide-issue superscalar processor. Higher performance is enabled thanks to VP (see Fig. 12) while EOLE enables a much simpler and far less power hungry OoO engine. The extra hardware blocks required for EOLE are relatively simple: Sets of ALUs in Early Execution and Late Execution

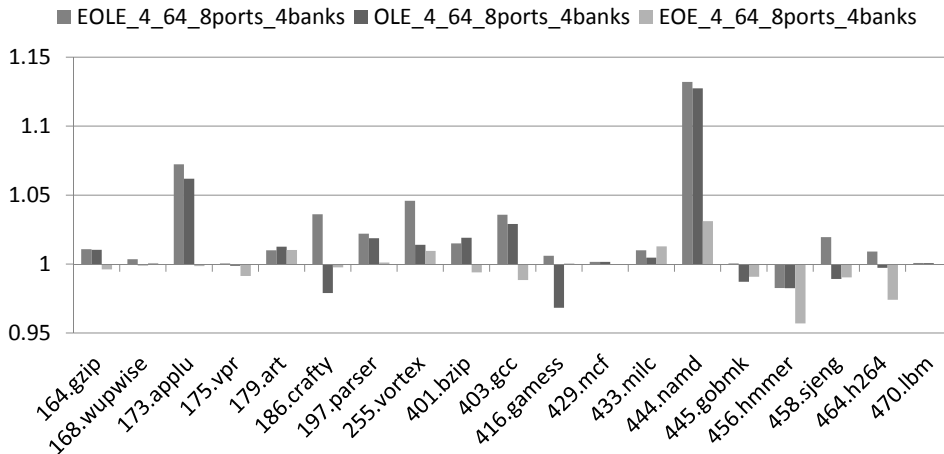


Figure 13: Performance of *EOLE_4_64*, *OLE_4_64* and *EOE_4_64* using 8 ports for Late Execution and validation and having 4 64-register banks, normalized to *Baseline_VP_6_64*.

stages, and storage tables and update logic for the value predictor itself.

6.5 A Note on the Modularity of EOLE: Introducing OLE and EOE

EOLE need not be implemented as a whole. In particular, either Early Execution or Late Execution can be implemented, if the performance vs. complexity tradeoff is deemed worthy. In particular, removing Late Execution can further reduce the number of read ports required on the PRF. Removing Early Execution saves on complexity since there is no need for an 8-to-8 bypass network anymore.

Fig. 13 shows the respective speedups of *EOLE_4_64*, *OLE_4_64* (Late Execution only) and *EOE_4_64* (Early Execution only) over *Baseline_VP_6_64*. As in the previous paragraph, only 8 read ports are dedicated to Late Execution/validation and the PRF is 4-banked (64 registers in each bank). The baseline has a single 256-register bank and enough ports to avoid contention.

We observe that some benchmarks are more sensitive to the absence of Late Execution (e.g. *applu*, *bzip*, *gcc*, *namd*, *hmmmer* and *h264*) while some are more sensitive to the absence of Early Execution (e.g. *crafty* and *gamess*). Nonetheless, the performance impact of removing Late Execution appears as more important in the general case.

In all cases slowdown over *Baseline_VP_6_64* remains under 5%. This suggests that when considering an effective implementation of VP using EOLE, an additional degree of freedom exists as either only Early or Late Execution may be implemented.

7 Conclusion and Future Work

Single thread performance remains the driving force for the design of new high-performance cores. However, hardware complexity and power consumption are still major obstacles to the implementation of new architectural features.

Value Prediction (VP) is one of such features that has still not been implemented in real-world products due to those obstacles. Fortunately, a recent advance in research on Value Prediction partially addressed these issues [24]. In particular, it was shown that prediction validation can

be performed at commit time without sacrificing performance. This removes the hardware logic associated with VP from the OoO engine and greatly simplifies design, as it eliminates the burdens of validation at execution-time and *selective replay*.

Building on VP validation at commit time, we have proposed EOLE, an *{Early / Out-of-Order / Late}* Execution microarchitecture aiming at further reducing the hardware complexity and the power consumption of a VP-augmented superscalar processor. By doing so, we strongly reinforce the case for an effective VP implementation in real hardware.

With Early Execution, single-cycle instructions whose operands are immediate or predicted are computed in-order in the front-end and do not have to flow through the OoO engine. With Late Execution, predicted single-cycle instructions as well as very high confidence branches are computed in-order in a pre-commit stage. They also do not flow through the OoO engine. As a result, EOLE significantly reduces the number of instructions dispatched to the OoO engine.

Considering an 6-wide, 64-entry IQ processor augmented with VP and prediction validation at commit time as the baseline, EOLE allows to drastically reduce the overall complexity and power consumption of both the OoO engine and the PRF. EOLE achieves performance very close to the baseline using only a 4-issue, 48-entry IQ OoO engine. It achieves similar or higher performance when using a 4-issue, 64-entry IQ engine, with the exception of one benchmark, *hammer* (1.8% slowdown).

With EOLE, the overhead over an 6-wide, 64-entry IQ processor (without VP) essentially consists of relatively simple hardware components, the two set of ALUs in the Early and Late Execution, a bypass network and the value predictor tables and update logic. The need for additional ports on the PRF is also substantially lowered by the reduction in issue width and some PRF optimizations (e.g. banking). The PRF can also be distributed into a copy in the OoO engine and a copy only read by the Late Execution/validation stage. Consequently, EOLE results in a much less complex and power hungry OoO engine, while generally benefiting from higher performance thanks to Value Prediction. Moreover, we hinted that Late Execution and Early Execution can be implemented separately, with Late Execution appearing as slightly more cost-effective.

Further studies to evaluate the possible variations of EOLE designs may include the full range of hardware complexity mitigation techniques that were discussed in Section 6.3 for both Early and Late execution, and the exploration of other possible sources of Late Execution, e.g. indirect jumps, returns, but also store address computations. One can also explore the interactions between EOLE and previous propositions aiming at reducing the complexity of the OoO engine such as the *Multicluster* architecture [8] or register file-oriented optimizations [37]. Further research also includes the need to look for more storage-effective value prediction schemes as well as even more accurate predictors.

Acknowledgment

This work was partially supported by the European Research Council Advanced Grant DAL No. 267175

References

- [1] P. Ahuja, D. Clark, and A. Rogers, “The performance impact of incomplete bypassing in processor pipelines,” in *Proceedings of the International Symposium on Microarchitecture*, 1995, pp. 36–45.

-
- [2] T. M. Austin, "DIVA: a reliable substrate for deep submicron microarchitecture design," in *Proceedings of the International Symposium on Microarchitecture*, 1999, pp. 196–207.
- [3] N. Binkert, B. Beckmann, G. Black, S. K. Reinhardt, A. Saidi, A. Basu, J. Hestness, D. R. Hower, T. Krishna, S. Sardashti, R. Sen, K. Sewell, M. Shoaib, N. Vaish, M. D. Hill, and D. A. Wood, "The gem5 simulator," *SIGARCH Comput. Archit. News*, vol. 39, no. 2, pp. 1–7, Aug. 2011.
- [4] G. Z. Chrysos and J. S. Emer, "Memory dependence prediction using store sets," in *Proceedings of the International Symposium on Computer Architecture*, 1998, pp. 142–153.
- [5] R. Eickemeyer and S. Vassiliadis, "A load-instruction unit for pipelined processors," *IBM Journal of Research and Development*, vol. 37, no. 4, pp. 547–564, 1993.
- [6] D. Ernst and T. Austin, "Efficient dynamic scheduling through tag elimination," in *Proceedings of the International Symposium on Computer Architecture*, 2002, pp. 37–46.
- [7] B. Fahs, T. Rafacz, S. J. Patel, and S. S. Lumetta, "Continuous optimization," in *Proceedings of the International Symposium on Computer Architecture*, 2005, pp. 86–97.
- [8] K. I. Farkas, P. Chow, N. P. Jouppi, and Z. Vranesic, "The Multicluster architecture: reducing cycle time through partitioning," in *Proceedings of the International Symposium on Microarchitecture*, 1997, pp. 149–159.
- [9] B. Fields, S. Rubin, and R. Bodik, "Focusing processor policies via critical-path prediction," in *Proceedings of the International Symposium on Computer Architecture*, 2001, pp. 74–85.
- [10] F. Gabbay and A. Mendelson, "Using value prediction to increase the power of speculative execution hardware," *ACM Trans. Comput. Syst.*, vol. 16, no. 3, pp. 234–270, Aug. 1998.
- [11] S. Gochman, R. Ronen, I. Anati, A. Berkovits, T. Kurts, A. Naveh, A. Saeed, Z. Sperber, and R. C. Valentine, "The Intel Pentium M processor: Microarchitecture and performance," *Intel Technology Journal*, vol. 7, May 2003.
- [12] B. Goeman, H. Vandierendonck, and K. De Bosschere, "Differential FCM: Increasing value prediction accuracy by improving table usage efficiency," in *Proceedings of the International Conference on High-Performance Computer Architecture*, 2001, pp. 207–216.
- [13] Intel, *Intel 64 and IA-32 Architectures Software Developer's Manual*, May 2012. [Online]. Available: <http://www.intel.com/content/www/us/en/processors/architectures-software-developer-manuals.html>
- [14] S. Jourdan, R. Ronen, M. Bekerman, B. Shomar, and A. Yoaz, "A novel renaming scheme to exploit value temporal locality through physical register reuse and unification," in *Proceedings of the International Symposium on Microarchitecture*, 1998, pp. 216–225.
- [15] R. E. Kessler, E. J. Mclellan, and D. A. Webb, "The Alpha 21264 microprocessor architecture," in *Proceedings of the International Conference on Computer Design*, 1998, pp. 90–95.
- [16] I. Kim and M. H. Lipasti, "Half-price architecture," in *Proceedings of the International Symposium on Computer Architecture*, 2003, pp. 28–38.
- [17] —, "Understanding scheduling replay schemes," in *Proceedings of the International Symposium on High Performance Computer Architecture*, 2004, pp. 198–.

- [18] M. H. Lipasti and J. P. Shen, "Exceeding the dataflow limit via value prediction," in *Proceedings of the Annual International Symposium on Microarchitecture*. IEEE Computer Society, 1996, pp. 226–237.
- [19] M. Lipasti, C. Wilkerson, and J. Shen, "Value locality and load value prediction," *ASPLOS-VII*, 1996.
- [20] A. Lukefahr, S. Padmanabha, R. Das, F. Sleiman, R. Dreslinski, T. Wenisch, and S. Mahlke, "Composite cores: Pushing heterogeneity into a core," in *Proceedings of the International Symposium on Microarchitecture*, 2012, pp. 317–328.
- [21] A. Mendelson and F. Gabbay, "Speculative execution based on value prediction," Technion-Israel Institute of Technology, Tech. Rep. TR1080, 1997.
- [22] T. Nakra, R. Gupta, and M. Soffa, "Global context-based value prediction," in *Proceedings of the International Symposium On High-Performance Computer Architecture*, 1999, pp. 4–12.
- [23] S. Palacharla, N. Jouppi, and J. Smith, "Complexity-effective superscalar processors," in *Proceedings of the International Symposium on Computer Architecture*, 1997, pp. 206–218.
- [24] A. Perais and A. Seznec, "Practical data value speculation for future high-end processors," in *Proceedings of the International Symposium on High-Performance Computer Architecture*, 2014, to appear. [Online]. Available: <http://hal.inria.fr/hal-00904743>
- [25] E. Perelman, G. Hamerly, and B. Calder, "Picking statistically valid and early simulation points," in *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques*, 2003, pp. 244–.
- [26] V. Petric, T. Sha, and A. Roth, "Reno: a rename-based instruction optimizer," in *Proceedings of the International Symposium on Computer Architecture*, 2005, pp. 98–109.
- [27] B. Rychlik, J. Faistl, B. Krug, A. Kurland, J. Sung, M. Velev, and J. Shen, "Efficient and accurate value prediction using dynamic classification," *Carnegie Mellon University, CM μ ART-1998-01*, 1998.
- [28] Y. Sazeides and J. Smith, "The predictability of data values," in *Proceedings of the International Symposium on Microarchitecture*, 1997, pp. 248–258.
- [29] A. Seznec, "Storage free confidence estimation for the TAGE branch predictor," in *Proceedings of the International Symposium on High Performance Computer Architecture*, 2011, pp. 443–454.
- [30] A. Seznec and P. Michaud, "A case for (partially) TAgged GEometric history length branch prediction," *Journal of Instruction Level Parallelism*, vol. 8, pp. 1–23, 2006.
- [31] A. Seznec, E. Toullec, and O. Rochecouste, "Register write specialization register read specialization: a path to complexity-effective wide-issue superscalar processors," in *Proceedings of the International Symposium on Microarchitecture*, 2002, pp. 383–394.
- [32] Standard Performance Evaluation Corporation. CPU2000. [Online]. Available: <http://www.spec.org/cpu2000/>
- [33] ——. CPU2006. [Online]. Available: <http://www.spec.org/cpu2006/>

-
- [34] R. Thomas and M. Franklin, "Using dataflow based context for accurate value prediction," in *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques*, 2001, pp. 107–117.
 - [35] F. Tseng and Y. N. Patt, "Achieving out-of-order performance with almost in-order complexity," in *Proceedings of the International Symposium on Computer Architecture*, 2008, pp. 3–12.
 - [36] E. S. Tune, D. M. Tullsen, and B. Calder, "Quantifying instruction criticality," in *Proceedings of the Conference on Parallel Architectures and Compilation Techniques*, 2002, pp. 104–113.
 - [37] S. Wallace and N. Bagherzadeh, "A scalable register file architecture for dynamically scheduled processors," in *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques*, 1996, pp. 179–184.
 - [38] K. Wang and M. Franklin, "Highly accurate data value prediction using hybrid predictors," in *Proceedings of the International Symposium on Microarchitecture*, 1997, pp. 281–290.
 - [39] H. Zhou, J. Flanagan, and T. M. Conte, "Detecting global stride locality in value streams," in *In Proceedings of the Annual International Symposium on Computer Architecture*, 2003, pp. 324–335.
 - [40] V. Zyuban and P. Kogge, "The energy complexity of register files," in *Proceedings of the International Symposium on Low Power Electronics and Design*, 1998, pp. 305–310.



**RESEARCH CENTRE
RENNES – BRETAGNE ATLANTIQUE**

Campus universitaire de Beaulieu
35042 Rennes Cedex

Publisher
Inria
Domaine de Volveau - Rocquencourt
BP 105 - 78153 Le Chesnay Cedex
inria.fr

ISSN 0249-6399