

AN $\mathcal{O}(\log_2 N)$ -LATENCY SISO WITH APPLICATION TO BROADBAND TURBO DECODING

Peter A. Beerel and Keith M. Chugg
Electrical Engineering-Systems Dept.
University of Southern California
Los Angeles, CA 90089-2565

Abstract

The standard algorithm for computing the soft-inverse of a finite-state machine (i.e., the Soft-in/Soft-out or SISO) module, is the forward-backward algorithm. These forward and backward recursions can be computed in parallel, yielding an architecture with latency $\mathcal{O}(N)$, where N is the block size. We demonstrate that the standard SISO computation may be formulated using a combination of a prefix and suffix operations. Based on well-known tree-structures for fast parallel prefix computations in the Very Large Scale Integration literature (e.g., tree adders), we propose a tree-structured SISO that has latency $\mathcal{O}(\log_2 N)$. The decrease in latency comes primarily at a cost of area, with, in some cases, only a marginal increase in computation. We discuss how this structure could be used to design a very high throughput turbo decoder, or more generally an iterative detector. Various sub-windowing and tiling schemes are also considered to further improve latency.

I INTRODUCTION

Calculating the soft-inverse of a finite-state machine (FSM) is a key operation in many data detection/decoding algorithms. Perhaps the most appreciated application is iterative decoding of concatenated codes, such as turbo codes [1], [2]. However the SISO (soft-in/soft-out) module [3] is widely applicable in iterative and non iterative receivers and signal processing devices (e.g., [4], [5], [6], [7]). The soft-outputs generated by a SISO may also be thresholded to obtain optimal hard decision decoding (e.g., producing the same decisions as the Viterbi algorithm [8] or the Bahl algorithm [9]). The general trend in many applications is towards higher data rates and therefore fast algorithms and architectures are desired.

There are two performance (speed) aspects of a data detection circuit architecture that are relevant to this paper. The first is *throughput* which is a measurement of the number of bits per second the architecture can decode. The second is *latency* which is the end-to-end delay for decoding a block of N bits. *Non-pipelined* architectures are those that decode only one block at a time and for which the throughput is simply N divided by the latency. *Pipelined* architectures, on the other hand, may decode multiple blocks simultaneously, shifted in time [10], thereby achieving much higher throughput than their non-pipelined counterparts.

This work supported in part by the National Science Foundation (NCR-CCR-9726391).

Depending on the application, the throughput and/or latency of the data detection circuitry hardware is important. For example, the latency associated with interleaving in a turbo-coded system with relatively low data rate (less than 100Kb/s) will likely dominate the latency the iterative decoding hardware. For future high-rate systems, however, the latency due to the interleaver may become relatively small, making the latency of the decoder significant. While pipelined decoders [10] can often achieve the throughput requirements, such techniques generally do not substantially reduce latency. In addition, sometimes latency has a dramatic impact on overall system performance. For example, in a data storage system (e.g., magnetic hard drives), latency in the retrieval process has a dramatic impact on the performance of the microprocessor and the overall computer. Such magnetic storage channels use high-speed Viterbi processing with turbo-coded approaches suggested recently [11], [12].

The standard SISO algorithm is the forward-backward algorithm. The associated forward and backward recursions for all states can be computed in parallel, yielding an architecture with $\mathcal{O}(N)$ computational complexity and latency, where N is the block size. The key result of this paper is the re-formulation of the standard SISO computation using a combination of prefix and suffix operations, which leads to an architecture with $\mathcal{O}(\lg N)$ latency¹. This architecture is based on well-known tree-structures for fast parallel prefix computations in the Very Large Scale Integration (VLSI) literature (e.g., fast adders [13], [14]), so we refer to it as a *tree-SISO*.

This exponential decrease in latency for the tree-SISO comes at the expense of increased computational complexity and area. The exact value of these costs depends on the FSM structure, specifically the number of states, and the details of the implementation. However, for a four-state convolutional code, such as those often used as constituent codes in turbo codes, the tree-SISO architecture achieves $\mathcal{O}(\lg N)$ latency with computational complexity of $\mathcal{O}(N \lg N)$. Note that, for this four-state example, the computation complexity of tree-SISO architecture increases sublinearly with respect to the associated speed-up. This is better than well-studied linear-scale solutions to the Viterbi algorithm (e.g., [15]); the generalization of which to the SISO problem is not always clear. For this 4-state code example, the area associated with the $\mathcal{O}(\lg N)$ -tree-SISO is $\mathcal{O}(N)$.

After formally defining the SISO and prefix-suffix op-

¹We use \lg to denote \log_2 .

erations in Section II, we describe the reformulation and corresponding tree-SISO architecture in Sections III and IV, respectively. Compatibly of the tree-SISO with known latency reduction methods is discussed in Section V. We conclude with a discussion of the architecture's potential applications, feasibility, and performance given current VLSI trends.

II BACKGROUND

II-A Soft-In Soft-Out Modules

For concreteness, we consider a specific class of finite state machines with no parallel state transitions and a generic S -state trellis. Such a trellis has up to S transitions departing and entering each state. The FSM is defined by the labeling of the state transitions by the corresponding FSM input and FSM output. Let $t_k = (s_k, a_k, s_{k+1}) = (s_k, a_k) = (a_k, s_{k+1}) = (s_k, s_{k+1})$ be a trellis transition from state s_k at time k to state s_{k+1} in response to input a_k . Since there are no parallel state transitions, t_k is uniquely defined by any of these representations. Given that the transition t_k occurs, the FSM output is $x_k(t_k)$.²

Consider the FSM as a system that maps a digital input sequence a_k to a digital output sequence x_k . A marginal soft-inverse, or SISO, of this FSM can be defined as a mapping of soft-in (SI) information on the inputs $SI(a_k)$ and outputs $SI(x_k)$, to soft-output (SO) information for a_k and/or x_k . The mapping is defined by the combining and marginalization operators used. It is now well-understood that one need only consider one specific reasonable choice for marginalization and combining operators and the results easily translate to other operators of interest [13, Section 26.4], [16], [17]. Thus, we focus on the min-sum marginalization-combining operation with the results translated to max-product, sum-product, min*-sum, and max*-sum [18] in the standard fashion. In all cases, two functions K_1 and K_2 define the boundary time indices of a *combining window* or span used in the mapping for a particular quantity u_k (e.g., $u_k = s_k$, $u_k = a_k$, $u_k = t_k$, $u_k = x_k$, $u_k = (s_k, s_{k+d})$, etc.). For min-sum marginalization-combining, the *minimum sequence metric (MSM)* of a quantity u_k is the metric (or length) of the shortest path or sequence in a combining window or span that is consistent with the conditional value of u_k . Specifically, the MSM is defined as³

$$\text{MSM}_{K_1(k)}^{K_2(k)}(u_k) \triangleq \min_{\mathbf{t}_{K_1(k)}^{K_2(k)}:u_k} M_{K_1(k)}^{K_2(k)}(\mathbf{t}_{K_1(k)}^{K_2(k)}), \quad (1)$$

$$M_{K_1(k)}^{K_2(k)}(\mathbf{t}_{K_1(k)}^{K_2(k)}) \triangleq \sum_{m=K_1(k)}^{K_2(k)} M_m(t_m) \quad (2)$$

$$M(t_m) \triangleq \text{SI}(a_m) + \text{SI}(x_m(t_m)), \quad (3)$$

where the set of transitions starting at time K_1 and

²Note that for generality we allow the mapping from transitions to outputs to be dependent on k .

³We use the standard convention of an infinite-valued metric for transitions that cannot occur under the FSM structure.

ending at time K_2 that are consistent with u_k is denoted $\mathbf{t}_{K_1(k)}^{K_2(k)} : u_k$ and $\mathbf{t}_{K_1(k)}^{K_2(k)}$ implicitly defines a sequence of transitions $t_{K_1(k)}, t_{K_1(k)+1}, \dots, t_{K_2(k)}$. Depending on the specific application, one or both of the following will be computed

$$\text{SO}_{K_1(k)}^{K_2(k)}(x_k) \triangleq \text{MSM}_{K_1(k)}^{K_2(k)}(x_k) - \text{SI}(x_k) \quad (4)$$

$$\text{SO}_{K_1(k)}^{K_2(k)}(a_k) \triangleq \text{MSM}_{K_1(k)}^{K_2(k)}(a_k) - \text{SI}(a_k) \quad (5)$$

Because the system on which the SISO is defined is an FSM, the combining and marginalization operations in (2)-(3) can be computed efficiently. The traditional approach is the forward-backward algorithm which computes the MSM of the states recursively forward and backward in time. Specifically, for the standard fixed-interval algorithm based on soft-in for transitions t_k , $k = 0, 1, \dots, N-1$, we have the following recursion based on add-compare-select (ACS) operations

$$f_k(s_{k+1}) \triangleq \text{MSM}_0^k(s_{k+1}) \quad (6)$$

$$= \min_{t_k:s_{k+1}} [f_{k-1}(s_k) + M_k(t_k)] \quad (7)$$

$$b_k(s_k) \triangleq \text{MSM}_k^{N-1}(s_k) \quad (8)$$

$$= \min_{t_k:s_k} [b_{k+1}(s_{k+1}) + M_k(t_k)] \quad (9)$$

where $f_{-1}(s_0)$ and $b_N(s_N)$ are initialized according to available edge information. Note that, since there are S possible values for the state, these state metrics can be viewed as $(S \times 1)$ vectors \mathbf{f}_k and \mathbf{b}_k . The final soft-outputs in (4)-(5) are obtained by marginalizing over the MSM of the transitions t_k

$$\text{SO}_0^{N-1}(u_k) = \min_{t_k:u_k} [f_{k-1}(s_k) + M_k(t_k) + b_{k+1}(s_{k+1})] - \text{SI}(u_k) \quad (10)$$

where u_k is either x_k or a_k . We refer to the operation in (10) as a *completion operation*.

While the forward-backward algorithm is computationally efficient, straightforward implementations of it have large latency (i.e., $\mathcal{O}(N)$) due to ACS bottleneck in computing the causal and antcausal state MSM's.

II-B Prefix and Suffix Operations

A prefix operation is defined as a generic form of computation that takes in n inputs y_0, y_1, \dots, y_{n-1} and produces n outputs z_0, z_1, \dots, z_{n-1} according to the following [13, Section 29.2.2][14]:

$$z_0 = y_0 \quad (11)$$

$$z_i = y_0 \otimes \dots \otimes y_i, \quad (12)$$

where \otimes is any associative binary operator.

Similarly, a suffix operation can be defined as a generic form of computation that takes in n inputs y_0, y_1, \dots, y_{n-1} and produces n outputs z_0, z_1, \dots, z_{n-1} according to

$$z_{n-1} = y_{n-1} \quad (13)$$

$$z_i = y_i \otimes \dots \otimes y_{n-1}, \quad (14)$$

where \otimes is *any* associative binary operator. Notice that a suffix operation is simply a prefix operation anchored at the other edge.

Prefix and suffix operations are important since enable a class of efficient algorithms based on tree-structured architectures with low latency implementations. The most notable of which are the VLSI n -bit tree adders whose latency is $\mathcal{O}(\lg n)$ [13], [19], [14].

III REFORMULATION OF THE SISO OPERATION

The proposed low-latency architecture is derived by formulating the SISO computations in terms of a combination of a prefix and suffix operations. To obtain this formulation, define $C(s_k, s_m)$, for $m > k$, as the MSM of state pairs s_k and s_m based on the soft-inputs between them, i.e., $C(s_k, s_m) = \text{MSM}_k^{m-1}(s_k, s_m)$. The set of MSMs $C(s_k, s_m)$ can be considered an $(S \times S)$ matrix $\mathbf{C}(k, m)$. The causal state MSMs \mathbf{f}_{k-1} can be obtained from $\mathbf{C}(0, k)$ by marginalizing (e.g., minimizing) out over the condition on s_0 . The backward state metrics can be obtained in a similar fashion. Specifically,

$$f_{k-1}(s_k) = \min_{s_0} C(s_0, s_k) \quad (15)$$

$$b_k(s_k) = \min_{s_N} C(s_k, s_N) \quad (16)$$

With this observation, the key step of the algorithm is to compute $\mathbf{C}(0, k)$ and $\mathbf{C}(k, N)$ for $k = 0, 1, \dots, N-1$. Note that the inputs of the algorithm are the one-step transition metrics which can be written as $\mathbf{C}(k, k+1)$ for $k = 0, 1, \dots, N-1$. To show how this algorithm can be implemented with a prefix and suffix computation, we define a min-sum fusion operator on \mathbf{C} matrices that inputs two such matrices, one with a left-edge coinciding with the right-edge of the other, and marginalizes out the midpoint to obtain a pairwise state-MSM with larger span. Specifically, given $\mathbf{C}(k_0, m)$ and $\mathbf{C}(m, k_1)$, we define a \mathbf{C} Fusion Operator, or \otimes_C operator by

$$C(s_{k_0}, s_{k_1}) = C(s_{k_0}, s_m) \otimes_C C(s_m, s_{k_1}) \triangleq \min_{s_m} [C(s_{k_0}, s_m) + C(s_m, s_{k_1})]. \quad (17)$$

Note that the \otimes_C operator is an associative binary operator that accepts two matrices and returns one matrix. With this definition $\mathbf{C}(0, k)$ and $\mathbf{C}(k, N)$ for $k = 0, 1, \dots, N-1$ can be computed using the prefix and suffix operations as follows:

$$\begin{aligned} \mathbf{C}(0, k) &= \mathbf{C}(0, 1) \otimes_C \mathbf{C}(1, 2) \cdots \otimes_C \mathbf{C}(k-1, k) \\ \mathbf{C}(k, N) &= \mathbf{C}(k, k+1) \otimes_C \\ &\quad \cdots \otimes_C \mathbf{C}(N-2, N-1) \otimes_C \mathbf{C}(N-1, N). \end{aligned}$$

In general, a SISO algorithm can be based on the decoupling property of state-conditioning. Specifically, conditioning on all possible FSM state values at time k , the MSM shortest path problems on either side of this state condition may be solved independently and then fused together (e.g., as performed by the \mathbf{C} -fusion operator). More generally, the SISO operation can be

decoupled based on a partition of the observation interval with each subinterval processed independently and then fused together. For example, the forward-backward algorithm is based on a partition to the single-transition level with the fusing taking place *sequentially* in the forward and backward directions. In contrast, other SISO algorithms may be defined by specifying the partition and a schedule for fusing together the solutions to the sub-problems. This may be viewed as specifying an association scheme to the above prefix-suffix operations (i.e., grouping with parentheses). The \mathbf{C} -fusion operations may be simplified in some cases depending on the association scheme. For example, the forward-backward algorithm replaces all \mathbf{C} -fusion operations by the much simpler forward and backward ACSs. However, latency is also a function of the association scheme. In the next section, we present an architecture based on a pairwise tree-structured grouping. This structure allows only a small subset of the \mathbf{C} -fusion operations to be simplified, but facilitates a significant reduction in latency compared to the forward-backward algorithm, by fusing solutions to the subproblems in a *parallel*, instead of sequential, manner.

IV LOW-LATENCY TREE-SISO ARCHITECTURES

There are many known low-latency parallel architectures based on binary tree-structured groupings of prefix operations [19], [13], [14] that can be adopted to SISOs. All of these have targeted n -bit adder design where the binary associative operator is a simple 1-bit addition. In fact, to the best of our knowledge, this is the first application of parallel prefix-suffix architectures to an algorithm based on binary associative operators that are substantially more complex than 1-bit addition. The known parallel prefix architectures trade reduced area for higher latency and account for a secondary restriction of limited fanout of each computational module. This latter restriction is important when the computational modules are small and have delay comparable to the delay of wires and buffers (e.g., in adder design). The fusion operators, however, are relatively large. Consequently, given current VLSI trends, they will dominate the overall delay for the foreseeable future. Thus, we propose to adopt an architecture which minimizes latency with the minimal number of computational modules without regard to fanout [14].

Specifically, the forward and backward metrics, \mathbf{f}_{k-1} and \mathbf{b}_{N-k} , for $k = 1, 2, \dots, N$ can be obtained using a hierarchal tree-structure based on the *fusion-module* (FM) array shown in Fig. 1. We define a *complete set* of \mathbf{C} matrices on the interval $\{k_0, k_0 + K\}$ as the $2K-1$ matrices $\mathbf{C}(k_0, k_0 + m)$ and $\mathbf{C}(k_0 + m, k_0 + K)$ for $m = 1, 2, \dots, K-1$ along with $\mathbf{C}(k_0, k_0 + K)$. This is the MSM information for all state pairs on the span of K steps in the trellis with one state being either on the left or right edge of the interval. The module in Fig. 1 fuses the complete sets of \mathbf{C} matrices for two adjacent span- K intervals to produce a complete set of \mathbf{C} matrices on the combined span of size $2K$. Of the $4K-1$ output \mathbf{C} matrices, $2K$ are obtained from the $2(2K-1)$ inputs without any processing. The other $2K-1$ output \mathbf{C} matrices are obtained by $2K-1$ \mathbf{C} Fusion Modules,

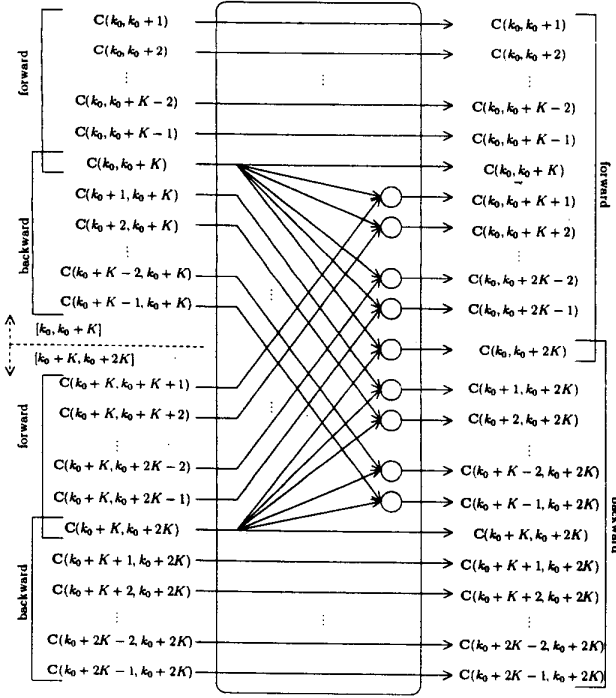


Fig. 1. The fusion processor array for combining the complete set of C matrices on $[k_0, k_0 + K]$ and $[k_0 + K, k_0 + 2K]$ to obtain the complete set on $[k_0, k_0 + 2K]$.

or CFMs, which implement the \otimes_C operator.

The basic span- K to span- $2K$ FM array shown in Fig. 1 can be utilized to compute the C matrices on the entire interval in $\lg N$ stages. This is illustrated in Fig. 2 for the special case of $N = 16$. Note that, indexing the stages from left to right (*i.e.*, increasing span) as $i = 1, 2, \dots, n = \lg N$ it is clear that there are 2^{n-i} FM arrays in stage i .

Because the final objective is to compute the causal and anticausal state metrics, however, not all FM's need be CFMs for all FM arrays. Specifically, the forward state metrics f_{k-1} can be obtained from f_{m-1} and $C(m, k)$ via

$$f_{k-1}(s_k) = \min_{s_m} [f_{m-1}(s_m) + C(s_m, s_k)] \quad (18)$$

Similarly, the backward state metrics can be updated via

$$b_k(s_k) = \min_{s_m} [b_m(s_m) + C(s_k, s_m)] \quad (19)$$

We refer to a processing module that produces an f vector from another f vector and a C matrix, as described in (18), as an f Fusion Module (*fFM*). A b Fusion Module (*bFM*) is defined analogously according to the operation in (19). In Fig. 2 we have indicated which FM's may be implemented as *f*FMs or *b*FMs.

The importance of this development is that the calculation of the state metrics has $\mathcal{O}(\lg N)$ latency. This is

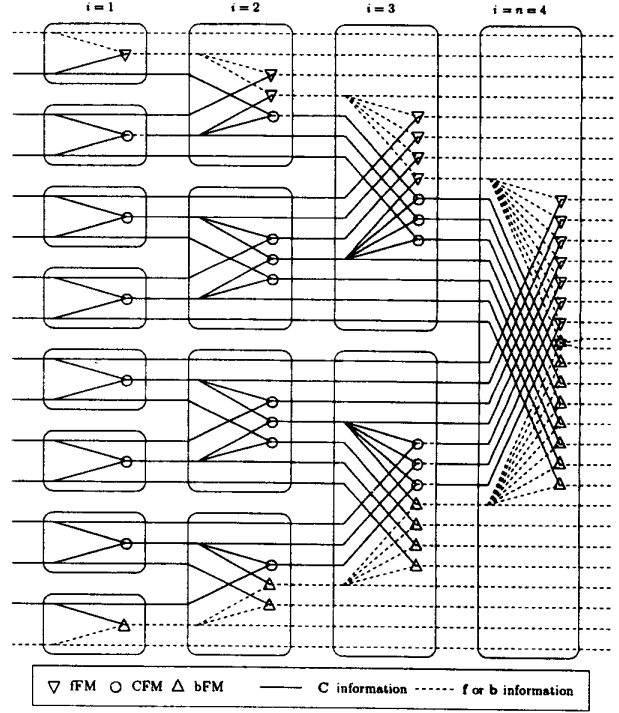


Fig. 2. The tree-SISO for $N = 16$.

because the only data dependencies are from one stage to the next and thus all FM arrays within a stage and all FMs within an FM array can be executed in parallel, each taking $\mathcal{O}(1)$ latency. The cost of this low latency is the need for relatively large amounts of area. One mitigating factor is that, because the stages of the tree operate in sequence, hardware can be shared between stages. Thus, the stage that requires the most hardware dictates the total hardware needed. A rough estimate of this is N sets of S S -way ACS units with the associated registers. For the example in Fig. 2, stage 2 has the most CFAs (8), but $i^* = 3$ has the most processing complexity. The complexity of stages $i = 1, 2, 3, 4$ is 26, 36, 32, and 16, respectively. Thus, if hardware is shared between stages, a total of 36 sets of S S -way ACS units is required to execute all FMs in a given stage in parallel. For applications when this number of ACS units is prohibitive, one can reduce the hardware requirements by as much as a factor of S with a corresponding linear increase in latency.

The implementation of the completion operation defined in (10) should also be considered. The basic operation required is a Q -way ACS unit where Q is the number of transitions consistent with u_k . Assuming that at most half of the transitions will be consistent with u_k , Q is upper bounded by $S^2/2$. Consequently, when S is large, low-latency, area-efficient implementations of the completion step may become an important issue. Fortunately, numerous low-latency implementations are

well-known [20]. The most straightforward may be one which uses a binary tree of comparators and has latency of $\mathcal{O}(\lg S^2)$. For small S , this additional latency is not significant.

The computational complexity of the state metric calculations can be computed using simple expressions based on Figs. 1 and 2. The total number of computations measured in units of S S -way ACS computations can be shown to be

$$N_{S,S} = N((\lg N - 3)S + 2) + 4S - 2 \quad (20)$$

For the example in Fig. 2, an equivalent of 110 sets of S S -way ACS operations are performed. This is to be compared with the corresponding forward-backward algorithm which would perform $2N = 32$ such operations and have baseline architectures with four times the latency. In general, note that for a reduction in latency from N to $\lg N$, the computation is increased by a factor of roughly $(1/2)(\lg N - 3)S + 1$. Thus, while the associated complexity is high, the complexity scaling is sub-linear in N .

IV-A Optimizations for Sparse Trellises

The above architecture is most efficient for fully-connected trellises. For sparser trellis structures, however, the initial processing modules must process C -matrices containing elements set to ∞ , accounting for MSMs of pairs of states between which there is no sequence of transitions, thereby wasting processing power and latency. This section discusses optimizations that address this inefficiency.

For concreteness, we consider as a baseline a standard 1-step trellis with M^L states and exactly M transitions into and out of each state, in which, there exists exactly one sequence of transitions to go from a given state at time s_k to a given state s_{k+L} . One optimization is to pre-collapse the one-step trellis into an R -step trellis, $1 \leq R \leq L$, and apply the tree SISO architecture to the collapsed trellis. A second optimization is to, wherever possible, simplify the C fusion modules. In particular, for a SISO on an R -step trellis, the first $\lg(L/R)$ stages can be simplified to banks of additions that simply add incoming pairs of multi-step transition metrics.

More precisely, pre-collapsing involves adding the R metrics of the 1-step transitions that constitute the transition metrics of each *super-transition* $\mathbf{t}_{kR}^{(k+1)R}$, for $k = 0, 1, \dots, (N-1)/R$. The SISO accepts these inputs and produces forward and backward MSMs, $f_{kR-1}(s_k)$ and $b_{(k+1)R}(s_{(k+1)R})$, for $k = 0, 1, \dots, N/R$. The key benefit of pre-collapsing is that number of SISO inputs is reduced by a factor of R , thereby reducing the number of stages required in the state metric computation by $\lg R$. One disadvantage of pre-collapsing is that the desired soft-outputs must be computed using a more complex, generalized completion operation. Namely,

$$\begin{aligned} \text{SO}_0^{NR-1}(u_{kR+m}) &= \min_{\mathbf{t}_{kL}^{(k+1)R}:u_{kR+m}} [f_{kR-1}(s_k) + \\ &M_{kR}^{(k+1)R}(\mathbf{t}_{kR}^{(k+1)R}) + b_{(k+1)R+1}(s_{(k+1)R})] \\ &- \text{SI}(u_{kR+m}) \quad m = 0, 1, \dots, R-1. \end{aligned} \quad (21)$$

The principle issue is that for each u_{kR+m} this completion step involves an $(M^{L+R}/2)$ -way ACS rather than the $(M^{L+1}/2)$ -way ACS required for the 1-step trellis.

In order to identify the optimal R assuming both these optimizations are performed, the relative latencies of the constitute operations are needed. While exact latencies are dependent on implementation details, rough estimates may still yield insightful results. In particular, we can assume that both the pre-collapsing additions and ACS operations for the state metric and completion operations are implemented using binary trees of adders/comparators and therefore estimate that their delay is logarithmic in the number of their inputs. An important observation is that the pre-collapsing along with $\lg R$ simplified stages together add L 1-step transition metrics (producing the transition metrics for a fully-connected L -step trellis) and thus can jointly be implemented in an estimated $\lg L$ time units. In addition, the state metric (M^L) -way ACS units take $\lg M^L$ time units and the completion units $(M^{L+R}/2)$ -way ACSs take $\lg(M^{L+R}/2)$ time units. Assuming maximal parallelism, this yields a total latency of

$$\lg L + \lg(N/R) \lg(M^L) + \lg(M^{L+R}/2) \quad (22)$$

It follows that the minimum latency occurs when $R - L \lg R$ is minimum (subject to $1 \leq R \leq L$), which occurs when $R = L$. This suggests that the minimum-latency architecture is one in which the trellis is pre-collapsed into a fully-connected trellis and more complex completion units are used to extract the soft outputs from the periodic state metrics calculated.

The cost of this reduced latency is the additional area required to implement the trees of adders that produce the L -step transition metrics and the larger trees of comparators required to implement the more complex completion operations. Note, however, that this area overhead can be mitigated by sharing adders and comparators among stages of each tree and, in some cases, between trees with only marginal impact on latency.

V USE IN TILED SUB-WINDOW SCHEMES

One known method of reducing latency and improving throughput of computing the soft-inverse is to use smaller combining windows. We define *minimum half-window (MHW)* algorithms as a class of SISOs in which the combining window bounds K_1 and K_2 satisfy $K_1(k) \leq \max(0, k-d)$ and $K_2(k) \geq \min(N, k+d)$, for $k = 0, \dots, N-1$ - i.e., for every point k away from the edge of the observation window, the the soft-output is based on a sub-window with left and right edges at least d points from k .

The traditional forward-backward algorithm can be used on sub-windows to obtain a MHW-SISO. One particular scheme is the *tiled sub-window technique* in which combining windows of length $2d+h$ are used to derive all state metrics. In this scheme, the windows are tiled with overlap of length $d+h$ and there are $\frac{N-2d}{h}$ such windows. Each sub-window yields h soft outputs, so there is an *overlap penalty* which increases as h decreases. For the

i^{th} such window, the forward and backward state metrics are computed using the recursions, modified from that of (7) and (9):

$$f_k^{(i)}(s_{k+1}) \triangleq \text{MSM}_{i(d+h)}^k(s_{k+1}) \quad (23)$$

$$b_k^{(i)}(s_k) \triangleq \text{MSM}_k^{d+(i+1)(d+h)-1}(s_k) \quad (24)$$

If all windows are processed in parallel, this yields a latency of $\mathcal{O}(2d+h)$.

The tree-SISO algorithm can be used in a MHW scheme without any overlap penalty and with $\mathcal{O}(\lg d)$ latency. Consider N/d combining windows of size d and let the tree-SISO compute $C(id, id+j)$ and $C((i+1)d, (i+1)d-j)$ for $j = 0, \dots, d-1$ and $i = 0, \dots, N/d-1$. Then, use one additional stage of logic to compute the forward and backward state metrics for all k time indices that fall within the i^{th} window, $i = 0, \dots, N/d-1$, as follows:⁴

$$\begin{aligned} f_k^{(i)}(s_{k+1}) &\triangleq \text{MSM}_{(i-1)d}^k(s_{k+1}) \\ &= \min_{s_{id}} \left\{ \left[\min_{s_{(i-1)d}} C(s_{(i-1)d}, s_{id}) \right] + C(s_{id}, s_{k+1}) \right\} \\ b_k^{(i)}(s_k) &\triangleq \text{MSM}_k^{(i+1)d}(s_k) \\ &= \min_{s_{id}} \left\{ C(s_k, s_{id}) + \left[\min_{s_{(i+1)d}} C(s_{id}, s_{(i+1)d}) \right] \right\} \end{aligned}$$

The inner minimization corresponds to a conversion from **C** information to **f** (b) information as in (15)-(16). The outer minimization corresponds to an fFM or bFM. The order of this minimization was chosen to minimize complexity. This is reflected in the example of this approach shown in Fig. 3, where the last stage of each of the four tree-SISOs is modified to execute the above minimizations in the proposed order. We refer to the module that does this as a *2Cfb module*. This may be viewed as a specialization of the stage 2 center CFMs in Fig. 2. The above combining of sub-window tree-SISO outputs adds one additional processing stage so that the required number of stages of FMs is $\lg(d) + 1$.

V-A Computational Complexity Comparison

The computational complexity of computing the state metrics using the forward-backward tiled scheme is the number of windows times the complexity of computing the forward-backward algorithm on each window. In terms of S S -way ACSs, this can be approximated for large N via

$$\frac{N-2d}{h} 2(d+h) \approx \frac{2N}{h} (d+h). \quad (25)$$

The computational complexity of computing the state metrics using the tree-SISO tiled scheme in terms of S S -way ACSs can be developed similarly and is

$$\frac{N}{d} d \lg(d) S + 2N = N(S \lg(d) + 2). \quad (26)$$

⁴This should be interpreted with $C(s_{-d}, s_0)$ replaced by initial left-edge information and similarly for $C(s_{N-1}, s_{N+d-1})$.

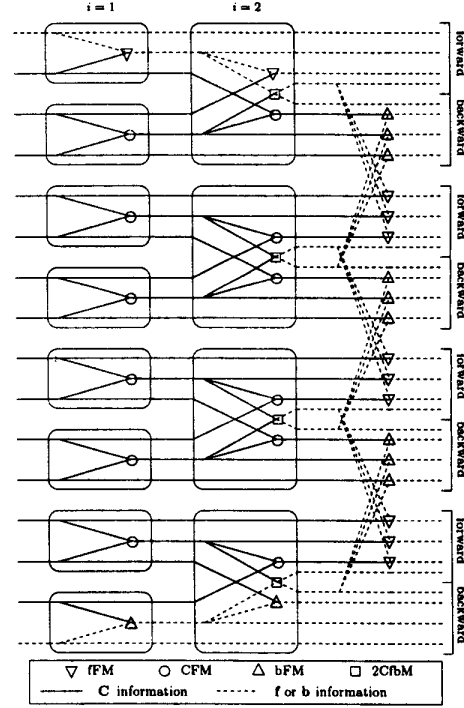


Fig. 3. A tiled sub-window approach with 4 tree-SISOs of window size 4 for $N = 16$.

Determining which scheme has higher computational complexity depends on the relative sizes of h and d . If h is reduced, the standard forward-backward scheme reduces in latency but increases in computational complexity because the number of overlapped windows increase. Since the tiled tree-SISO architecture has no overlap penalty, as h is decreased in a tiled forward-backward scheme, the relative computational complexity trade-off becomes more favorable for the tree-SISO approach. In fact, for $h < \frac{2d}{S \lg d}$, the computational complexities of the tree-SISO is lower than the tiled forward-backward scheme.

VI A DESIGN EXAMPLE: 4-STATE PCCC

The highly parallel architectures considered require large implementation area. In this section we consider an example for which the area requirements are most feasible for implementation in the near future. Specifically, we consider an iterative decoder based on 4-state sparse (one-step) trellises. Considering larger S will yield more impressive latency reductions for the tree-SISO. This is because the latency-reduction obtained by the tree-SISO architecture relative the parallel tiled forward backward architecture depends on the minimum half-window size. One expects that good performance requires a value of d that grows with the number of states (*i.e.*, similar to the rule-of-thumb for traceback depth in the Viterbi algorithm [21] for sparse trellises).

In contrast, considering pre-collapsing will yield less impressive latency reductions. For example, if $d = 16$ is required for a single-step trellis, then an effective value of $d = 8$ would suffice for a two-step trellis. The latency reduction factor associated with the tree-SISO for the former would be approximately 4, but only $8/3$ for the latter. However, larger S and/or pre-collapsing yields larger implementation area and is not in keeping with our desire to realistically assess the near-term feasibility of these algorithms.

In particular, we consider a standard parallel concatenated convolutional code (PCCC) with two 4-state constituent codes [1], [2]. Each of the recursive systematic constituent codes generates parity using the generator polynomial $G(D) = (1 + D^2)/(1 + D + D^2)$ with parity bits punctured to achieve an overall systematic code with rate $1/2$.

In order to determine the appropriate value for d to be used in the MHW-SISOs, we ran simulations where each SISO used a combining window $\{k - d, \dots, k + d\}$ to compute the soft-output at time k . This is exactly equivalent to the SISO operation obtained by a tiled forward-backward approach with $h = 1$. Note that, since d is the size of all (interior) half-windows for the simulations, any architecture based on a MHW-SISO with d will perform at least as well (e.g., $h = 2$ tiled forward-backward, d -tiled tree-SISO, etc.). For an interleaver size of $N = 1024$ with min-sum marginalization and combining and ten iterations, $d = 16$ yields performance similar to the fixed-interval curve [22]. This is consistent with the rule-of-thumb of five to seven times the memory for the traceback depth in a Viterbi decoder (i.e., roughly $d = 7 \times 2 = 14$ is expected to be sufficient).

Since the required window size is $d = 16$, the latency improvement of a tree-SISO relative to a tiled forward-backward scheme is close to $4 = 16/\lg(16)$. The computational complexity of these two approaches is similar and depends on the details of the implementation and the choice of h for the tiled forward-backward approach. A complete fair comparison would require a detailed implementation of the two approaches. Below we summarize a design for the tree-SISO based sub-window architecture.

A factor that impacts the area of the architecture is the bit-width of the data units. Simulation results suggest that an 8-bit datapath is sufficient. Roughly speaking, a tree-based architecture for this example would require 1024 sets of sixteen 4-way ACS units along with associated output registers to store intermediate state metric results. Each 4-way ACS unit can be implemented with an 8-bit 4-to-1 multiplexor, four 8-bit adders, six 8-bit comparators, and one 8-bit register [20]. Our initial VLSI designs indicate that these units require approximately 2250 transistors. Thus, this yields an estimate of $16 \times 2250 \times 1024 \approx 40$ Million transistors. This number or logic transistors pushes the limit of current VLSI technology but should soon be feasible. We consider an architecture in which one clock cycle is used per stage of the tree at a 200 MHz clock frequency. For $d = 16$, each SISO operation can be performed in 6 such clock cycles (using one clock for the completion step). Moreover, we assume a hard-wired interleaver compris-

ing two rows of 1024 registers with interconnection a network. Such an interleaver would be larger than existing memory-based solutions [10], but could have a latency of 1 clock cycle. Consequently, one iteration of the turbo decoder, consisting of two applications of the SISO, one interleaving, and one deinterleaving, requires an 14 clock cycles. Assuming ten iterations, the decoding of 1024 bits would take 140 clock cycles, or a latency of just 700 ns.

This latency also implies a very high throughput which can further be improved with standard pipelining techniques. In particular, a non-pipelined implementation has an estimated throughput of 1024 bits per 700 ns = 1.5 Gb/second. Using the tree-SISO architecture one could also pipeline across interleaver blocks as described by Masera et al. [10]. In particular, 20 such tiled tree-SISOs and associated interleavers can be used to achieve a factor of 20 in increased throughput, yielding a throughput of 30 Gb/second.

Moreover, unlike architectures based on the forward-backward algorithm, the tree-SISO can easily be internally pipelined, yielding even higher throughputs with linear hardware scaling. In particular, if dedicated hardware is used for each stage of the tree-SISO, pipelining the tree-SISO internally may yield another factor of $\lg(d)$ in throughput, with no increase in latency. For window sizes of $d = 16$, the tree-based architecture could support over 180 Gb/second. That said, it is important to realize that with current technology such hardware costs may be beyond practical limits. Given the continued increasing densities of VLSI technology, however, even such aggressive architectures may become cost-effective in the future.

VII CONCLUSION

Based on the interpretation of the SISO operation in terms of parallel prefix/suffix operations, a family of tree-structured architectures were suggested. Compared to the baseline forward-backward algorithm architecture, the tree-SISO architecture reduces latency from $\mathcal{O}(N)$ to $\mathcal{O}(\lg N)$. The tree-SISO was also demonstrated to be compatible with tiled sub-window approaches. Latency in this case is reduced from linear in the minimum half-window size (d) for fully-parallel tiled architectures based on the forward-backward algorithm, to logarithmic in d for tiled tree-SISOs.

The potential latency advantages clearly are most significant for applications requiring large combining windows. This is expected when the number of states increases. In the one detailed 4-state example considered, the latency was reduced by a factor of approximately 4. For systems with binary inputs and S states, one would expect that $d \cong 8 \lg(S)$ would be sufficient. Thus, there is a potential reduction in latency of approximately $8 \lg(S)/\lg(8 \lg S)$ which becomes quite significant as S increases. However, the major challenge in achieving this potential latency improvement is the area required for the implementation. In particular, building a high-speed S -way ACS unit for large S is the key challenge. In addition, alternative architectures that trade increase latency for reduced area requirements are also possible. In fact, we recently discovered an alternative tree ar-

chitecture that incurs a linear increase in latency with respect to the tree-SISO described herein, with significantly lower computational complexity and area requirements [23].

REFERENCES

- [1] C. Berrou, A. Glavieux, and P. Thitimajshima, "Near shannon limit error-correcting coding and decoding: turbo-codes," in *International Conference on Communications*, (Geneva, Switzerland), pp. 1064–1070, May 1993.
- [2] C. Berrou and A. Glavieux, "Near optimum error correcting coding and decoding: turbo-codes," *IEEE Trans. Commun.*, vol. 44, pp. 1261–1271, October 1996.
- [3] S. Benedetto, D. Divsalar, G. Montorsi, and F. Pollara, "Soft-input soft-output modules for the construction and distributed iterative decoding of code networks," *European Trans. Commun.*, vol. 9, March/April 1998.
- [4] X. Chen and K. M. Chugg, "Near-optimal page detection for two-dimensional ISI/AWGN channels using concatenated modeling and iterative detection," in *Proc. International Conf. Communications*, (Atlanta, GA), 1998.
- [5] K. M. Chugg, X. Chen, A. Ortega, and C.-W. Chang, "An iterative algorithm for two-dimensional digital least metric problems with applications to digital image compression," in *Proc. Intl. Conf. Image Processing*, (Chicago, IL), 1998.
- [6] M. Moher and P. Guinand, "An iterative algorithm for asynchronous coded multiuser detection," *IEEE Communications Letters*, August 1998.
- [7] X. Chen and K. M. Chugg, "Reduced state soft-in/soft-out for complexity reduction in iterative and non-iterative data detection," in *Proc. International Conf. Communications*, (New Orleans, LA), 2000.
- [8] G. D. Forney, "The Viterbi algorithm," *Proc. IEEE*, vol. 61, pp. 268–278, March 1973.
- [9] L. R. Bahl, J. Cocke, F. Jelinek, and J. Raviv, "Optimal decoding of linear codes for minimizing symbol error rate," *IEEE Trans. Inform. Theory*, vol. IT-20, pp. 284–287, March 1974.
- [10] G. Masera, G. Piccinini, M. R. Roch, and M. Zamboni, "VLSI architectures for turbo codes," *IEEE Transactions on VLSI*, vol. 7, September 1999.
- [11] A. Bishop, I. Chan, S. Aronson, P. Moran, K. Hen, R. Cehng, K. K. Fitzpatrick, J. Stander, R. Chik, K. Kshonze, M. Ali-ahmad, J. Ngai, H. He, E. daVeiga, P. Bolte, C. Krasuk, B. Cerqua, and R. Brown, "A 300Mb/s BiCMOS disk drive channel with adaptive analog equalizer," in *IEEE International Solid-State Circuits Conference*, 1999.
- [12] W. Ryan, "Performance of high rate turbo codes on a pr4-equalized magnetic recording channel," in *Proc. International Conf. Communications*, (Atlanta, GA), 1998.
- [13] T. H. Cormen, C. Leiserson, and R. L. Rivest, *Introduction to Algorithms*. Cambridge, Mass.: The MIT Press, 1990.
- [14] A. E. Despain, "Chapter 2: Notes on computer architecture for high performance," in *New Computer Architectures*, Academic Press, 1984.
- [15] G. Fettweis and H. Meyr, "Parallel Viterbi algorithm implementation: Breaking the ACS-bottleneck," *IEEE Trans. Commun.*, vol. 37, pp. 785–790, August 1989.
- [16] N. Wiberg, *Codes and Decoding on General Graphs*. PhD thesis, Linköping University (Sweden), 1996.
- [17] S. M. Aji and R. J. McEliece, "The generalized distributive law," *IEEE Trans. Inform. Theory*, vol. 46, pp. 325–343, March 2000.
- [18] P. Robertson, E. Vilebrum, and P. Hoeher, "A comparison of optimal and suboptimal MAP decoding algorithms operating in the log domain," in *Proc. International Conf. Communications*, (Seattle, WA), pp. 1009–1013, 1995.
- [19] R. P. Brent and H. T. Kung, "A regular layout for parallel adders," *IEEE Transactions on Computers*, vol. C-31, pp. 260–264, March 1982.
- [20] P. J. Black, *Algorithms and Architectures for High Speed Viterbi Decoding*. PhD thesis, Stanford University, California, March 1993.
- [21] J. A. Heller and I. Jacobs, "Viterbi decoding for satellite and space communication," *IEEE Trans. Commun. Technol.*, vol. COM-19, pp. 835–848, October 1971.
- [22] K. M. Chugg, A. Anastasopoulos, and X. Chen, *Iterative Detection: Adaptivity, Complexity Reduction, and Applications*. Kluwer Academic Publishers, 2000.
- [23] P. Thienviboon and K. M. Chugg, "A low-latency SISO via message passing on a binary tree," in *Allerton Conf. (submitted)*, 2000.