
Relaxed Scheduling for Scalable Belief Propagation

Abstract

1 The ability to leverage large-scale hardware parallelism has been one of the key
2 enablers of the accelerated recent progress in machine learning. Consequently,
3 there has been considerable effort invested into developing efficient parallel variants
4 of classic machine learning algorithms. However, despite the wealth of knowledge
5 on parallelization, some classic machine learning algorithms often prove hard to
6 parallelize efficiently while maintaining convergence.

7 In this paper, we focus on efficient parallel algorithms for the key machine learning
8 task of inference on graphical models, in particular on the fundamental belief
9 propagation algorithm. We address the challenge of efficiently parallelizing this
10 classic paradigm by showing how to leverage scalable relaxed schedulers in this
11 context. We present an extensive empirical study, showing that our approach
12 outperforms previous parallel belief propagation implementations both in terms
13 of scalability and in terms of wall-clock convergence time, on a range of practical
14 applications.

15 1 Introduction

16 *Hardware parallelism* has been a key computational enabler for recent advances in machine learning,
17 as it provides a way to reduce the processing time for the ever-increasing quantities of data required for
18 training accurate models. Consequently, there has been considerable effort invested into developing
19 efficient parallel variants of classic machine learning algorithms, e.g. [28, 22, 23, 24, 16].

20 In this paper, we will focus on efficient parallel algorithms for the fundamental task of *inference on*
21 *graphical models*. The inference task in graphical models takes the form of *marginalisation*: we are
22 given observations for a subset of the random variables, and the task is to compute the conditional
23 distribution of one or a few variables of interest. The marginalization problem is known to be
24 computationally intractable in general [10, 33, 9], but inexact heuristics are well-studied for practical
25 inference tasks.

26 One popular heuristic for inference on graphical models is *belief propagation* [27], inspired by the
27 exact dynamic programming algorithm for marginalization on trees. While belief propagation has
28 no general approximation or even convergence guarantees, it has proven empirically successful in
29 inference tasks, in particular in the context of decoding low-density parity check codes [8]. However,
30 it remains poorly understood how to properly parallelize belief propagation.

31 **Parallelizing belief Propagation.** To illustrate the challenges of parallelizing belief propagation,
32 we will next give a simplified overview of the belief propagation algorithm, and refer the reader
33 to Section 2 for full details. Belief propagation can be seen as a *message passing* or a *weight*
34 *update* algorithm. In brief, belief propagation operates over the underlying graph $G = (V, E)$ of the
35 graphical model, maintaining a vector of real numbers called a *message* $\mu_{i \rightarrow j}$ for each ordered pair
36 (i, j) corresponding to an edge $\{i, j\} \in E$ (Fig. 1). The core of the algorithm is the *message update*
37 *rule* which specifies how to update an outgoing message $\mu_{i \rightarrow j}$ at node i based on the *other* incoming
38 messages at node i ; for the purposes of the present discussion, it is sufficient to view this as black
39 box function f over these other messages, leading to the update rule

$$\mu_{i \rightarrow j} \leftarrow f(\{\mu_{k \rightarrow i} : k \in N(i) \setminus \{j\}\}). \quad (1)$$

40 This update rule is applied to messages until the values of messages have converged to a stable
41 solution, at which point the algorithm is said to have terminated.

42 Importantly, the message update rule does not
 43 specify *in which order* messages should be
 44 updated. The standard solution, called *syn-*
 45 *chronous belief propagation*, is to update all the
 46 message simultaneously. That is, in each global
 47 round $t = 1, 2, 3, \dots$, given message values
 48 $\mu_{i \rightarrow j}^t$ for all pairs (i, j) , the new values $\mu_{i \rightarrow j}^{t+1}$
 49 are computed as

$$\mu_{i \rightarrow j}^{t+1} \leftarrow f(\{\mu_{k \rightarrow i}^t : k \in N(i) \setminus \{j\}\})$$

50 However, there is evidence that updating mes-
 51 sages *one at a time* leads to faster and more
 52 reliable convergence [14]; in particular, various
 53 proposed *priority-based schedules*—schedules
 54 that try to prioritize message updates that would make ‘more progress’—have proven empirically to
 55 converge with much fewer message updates than the synchronous schedule [14, 20, 38].

56 Having to execute updates in a strict priority order poses a challenge for efficient *parallel* imple-
 57 mentations of belief propagation: while the synchronous schedule is naturally parallelizable, as all
 58 message updates can be done independently, the more efficient priority-based schedules are inherently
 59 sequential and thus seem difficult to parallelize. Accordingly, existing work on efficient parallel belief
 60 propagation has focused on designing custom schedules that try to import some features from the
 61 priority-based schedules while maintaining a degree of parallelism [16, 11].

62 **Our contributions.** In this work, we address the challenges of parallel belief propagation by
 63 showing how to efficiently parallelize any priority-based schedule for belief propagation. The key
 64 idea is that we can *relax* the priority-based schedules by allowing limited out-of-order execution,
 65 concretely implemented using a *relaxed scheduler*, as we will explain next.

66 Consider a belief propagation algorithm that schedules the message updates according to a priori-
 67 ty function r by always updating the message $\mu_{i \rightarrow j}$ with the highest priority $r(\mu_{i \rightarrow j})$ next; this
 68 framework captures existing priority-based schedules such as residual belief propagation [14] and its
 69 variants [20, 38]. Concretely, an iterative centralized version of this algorithm can be implemented
 70 by storing the messages in a priority queue Q , and iterating the following procedure:

- 71 (1) Pop the top element for Q to obtain the message $\mu_{i \rightarrow j}$ with highest priority $r(\mu_{i \rightarrow j})$.
- 72 (2) Update message $\mu_{i \rightarrow j}$ following (1).
- 73 (3) Update the priorities in Q for messages affected by the update.

74 This template does not easily lend itself to efficient parallelization, as the priority queue Q becomes a
 75 contention bottleneck. Previous work, e.g. [16, 11] investigated various heuristics for the parallel
 76 scheduling of updates in belief propagation, trading off increased parallelism with additional work in
 77 processing messages or even potential loss of convergence.

78 In this paper, we investigate an alternative approach, replacing the priority queue Q with a *relaxed*
 79 *priority queue (scheduler)* to obtain an efficient parallel version of the above template. A *relaxed*
 80 *scheduler* [3, 1, 2, 5] is similar to a priority queue, but instead of guaranteeing that the *top* element
 81 is always returned first, it only guarantees to return *one of the top k elements by priority*, where k
 82 is a variable parameter. Relaxed schedulers are popular in the context of parallel graph processing,
 83 e.g. [17, 26], and can induce non-trivial trade-offs between the degree of relaxation and the scalability
 84 of the underlying implementation, e.g. [1, 5].

85 For belief propagation, relaxed schedulers induce a *relaxed priority-based scheduling* of the mes-
 86 sages, roughly following the original schedule but allowing for message updates to be performed
 87 out of order, with guarantees on the maximum degree of priority inversion. We investigate the
 88 scalability-convergence trade-off between the *degree of relaxation* in the scheduler, and the *conver-*
 89 *gence behaviour* of the underlying algorithm, both theoretically and practically. In particular:

- 90 – We present a general scheme for parallelizing belief propagation schedules using relaxed sched-
 91 ulers with theoretical guarantees. While relaxed schedulers have been applied in other settings,
 92 and relaxed scheduling has been touched upon in belief propagation scheduling [16], no systematic
 93 study on relaxed belief propagation scheduling has been performed in prior work.

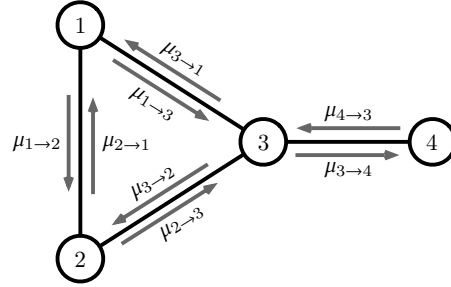


Figure 1: State of the belief propagation algorithm consist of two directed messages for each edge.

- 94 – We provide a theoretical analysis on the effects of relaxed scheduling for belief propagation on
- 95 trees. We exhibit both positive results—instance classes where relaxation overhead is negligible—
- 96 and negative results, i.e., worst-case instances where relaxation can cause significant wasted
- 97 work.
- 98 – We implement and experimentally compare different variants of belief propagation under relaxed
- 99 scheduling. We identify a new family of relaxed schedulers which consistently matches or
- 100 outperforms previous proposals. Benchmarks show that this framework gives state-of-the-art
- 101 parallel scalability on a wide variety of Markov random field models.

102 2 Preliminaries and related work

103 2.1 Belief Propagation

104 We consider marginalization in *pairwise Markov random fields*; one can equivalently consider factor
 105 graphs or Bayesian networks [40]. A pairwise Markov random field is defined by a set of random
 106 variables X_1, X_2, \dots, X_n , a graph $G = (V, E)$ with $V = \{1, 2, \dots, n\}$, and a set of *factors*

$$\begin{aligned} \psi_i &: D_i \rightarrow \mathbb{R}^+ && \text{for } i \in V, \\ \psi_{ij} &: D_i \times D_j \rightarrow \mathbb{R}^+ && \text{for } \{i, j\} \in E, \end{aligned}$$

107 where D_i denotes the domain of random variable X_i . The edge factors ψ_{ij} represent the dependencies
 108 between the random variables, and the node factors ψ_i represent a priori information about the
 109 individual random variables; the Markov random field defines a joint probability distribution on
 110 $X = (X_1, X_2, \dots, X_n)$ as

$$\Pr[X = x] \propto \prod_i \psi_i(x_i) \prod_{ij} \psi_{ij}(x_i, x_j),$$

111 where the ‘proportional to’ notation \propto hides the normalization constant applied to the right-hand
 112 side to obtain a probability distribution. The marginalization problem is to compute the probabilities
 113 $\Pr[X_i = x]$ for a specified subset of variables; for convenience, we assume that any observations
 114 regarding the values of other variables are encoded in the node factor functions ψ_i .

115 Belief propagation is a message-passing algorithm; for each ordered pair (i, j) such that $\{i, j\} \in E$,
 116 we maintain a *message* $\mu_{i \rightarrow j} : D_j \rightarrow \mathbb{R}$, and the algorithm iteratively updates these messages until
 117 the values (approximately) converge to a fixed point. On Markov random fields, the message update
 118 rule gives the new value of message $\mu_{i \rightarrow j}$ as a function of the old messages directed to node i by

$$\mu_{i \rightarrow j}(x_j) \propto \sum_{x_i \in D_i} \psi_i(x_i) \psi_{ij}(x_i, x_j) \prod_{k \in N(i) \setminus \{j\}} \mu_{k \rightarrow i}(x_i), \quad (2)$$

119 where $N(j)$ denotes the neighbors of node j in the graph G . Once the algorithm has converged, the
 120 marginals are estimated as

$$\Pr[X_i = x_i] \propto \psi_i(x_i) \prod_{j \in N(i)} \mu_{j \rightarrow i}(x_i).$$

121 The update rule (2) can be applied in arbitrary order. The standard *synchronous belief propagation*
 122 updates all the message simultaneously; in each global round $t = 1, 2, 3, \dots$, given message values
 123 $\mu_{i \rightarrow j}^t$ for all pairs $\{i, j\} \in E$, the new values $\mu_{i \rightarrow j}^{t+1}$ are computed as

$$\mu_{i \rightarrow j}^{t+1}(x_j) \propto \sum_{x_i \in D_i} \psi_i(x_i) \psi_{ij}(x_i, x_j) \prod_{k \in N(i) \setminus \{j\}} \mu_{k \rightarrow i}^t(x_i).$$

124 **Asynchronous belief propagation.** Starting with Elidan et al. [14], there has been a line of research
 125 arguing that *asynchronous* or *iterative* schedules for belief propagation tend to converge more reliably
 126 and with fewer message updates than the synchronous schedule. In particular, practical work has
 127 focused on schedules that attempt to iteratively perform ‘the most useful’ update at each step; the
 128 most prominent of these algorithms is the *residual belief propagation* of Elidan et al. [14], with other
 129 proposals aiming to address its shortcomings in various cases.

130 **Residual belief propagation.** Given a current state of messages, let $\mu'_{i \rightarrow j}$ denote the message we
 131 would obtain by applying the message update rule (2) to message $\mu_{i \rightarrow j}$. In residual belief propagation,
 132 the priority of a message is given by the *residual* $\text{res}(\mu_{i \rightarrow j})$ of a message $\mu_{i \rightarrow j}$, defined as

$$\text{res}(\mu_{i \rightarrow j}) = \|\mu'_{i \rightarrow j} - \mu_{i \rightarrow j}\|, \quad (3)$$

133 where $\|\cdot\|$ is an arbitrary norm; in this work, we assume L^2 norm is used unless otherwise specified.
134 That is, the residual of a message corresponds to amount of change that would happen if message
135 $\mu_{i \rightarrow j}$ would be updated. Note that this means that residual belief propagation performs *lookahead*,
136 that is, the algorithm precomputes the future updates before applying them to the state of the algorithm.
137 We will explore the performance of this base algorithm, as well as additional variants with *weight*
138 *decay* [20] and *without lookahead* [38]. (Due to space constraints, we leave their detailed description
139 to Appendix C.)

140 2.2 Parallel belief propagation

141 The question of parallelizing belief propagation is not fully understood. The synchronous schedule
142 is trivially parallelizable by performing updates within each round in parallel, but the improved
143 convergence properties of iterative schedules cannot easily be translated to parallel setting. Recent
144 proposals aim to bridge this gap in an ad-hoc manner by designing custom algorithms for specific
145 parallel settings.

146 **Residual splash.** *Residual splash* [16] is a vertex-based algorithm inspired by residual BP. Residual
147 splash was designed for MapReduce computation, and it aims to have larger individual tasks while
148 retaining a similar structure to residual BP. Specifically, the algorithm works by defining a priority
149 function over nodes of the Markov random field, and selecting the next node to process in a strict
150 priority order. For the selected node, the algorithm performs a *splash* operation that propagates
151 information within distance H in the graph; in practice, this results in threads performing larger
152 individual tasks at once, offsetting the cost of accessing the strict scheduler.

153 **Randomized synchronous belief propagation.** Van der Merve et al. [11] proposed a parallelization
154 scheme for belief propagation on GPUs, mixing the structure of synchronous and residual belief
155 propagation. Their algorithm considers all messages at once in global rounds, and performs *filter-*
156 *and-select* steps. First, it filters out all messages whose residuals are below the convergence threshold.
157 Second, out of the remaining messages, select a p -fraction uniformly at random to update. The
158 fraction p is adjusted on the fly based on the convergence of the algorithm, preferring a low value
159 if the algorithm is converging slowly, and a high value if it is converging fast. This algorithm is
160 well-suited for GPUs, as the filter-and-select steps can be efficiently implemented. However, as shown
161 by our experimental study, this strategy is not efficient on CPUs, on real-world models. Conversely,
162 as discussed in [11], the dynamic priority-based strategy we propose would be hard to implement
163 efficiently on GPUs, due to its irregular structure.

164 3 Relaxed priority-based belief propagation

165 In this section, we describe our framework for parallelizing belief propagation schedules via relaxed
166 schedulers. The main idea of the framework follows the description given in the introduction;
167 however, we generalize it to capture schedules that do not use individual messages as elementary
168 tasks, e.g. residual path [16].

169 3.1 Relaxed scheduling for iterative algorithms

170 Relaxed schedulers are a basic tool to parallelize iterative algorithms, used in the context of large-scale
171 graph processing [17, 26, 7, 12, 13]. At a high level, such iterative algorithms can be split into *tasks*,
172 each corresponding to a local operation involving, e.g., a vertex and its edges. A standard example is
173 Dijkstra’s classic shortest-paths algorithm, where each task updates the distance between a vertex
174 and the source, as well as the distances of the vertex’s neighbours. In many algorithms, tasks have
175 a natural priority order—in Dijkstra’s, the top task corresponds to the vertex of minimum distance
176 from the source. Many graph algorithms share this structure [37, 2], and can be mapped onto the
177 priority-queue pattern described in the introduction. However, due to contention on the priority queue,
178 implementing this pattern in practice can negate the benefits of parallelism [26].

179 **Relaxed Scheduler Definition.** A natural idea is to downgrade the guarantees of the perfect priority
180 queue, to allow for more parallelism. Relaxed schedulers [1] formalize this relaxation as follows.
181 We are given a parameter k , the degree of relaxation of the scheduler. The k -relaxed scheduler is a
182 *sequential* object supporting `Insert (<key, priority>)`, `IncreaseKey (<key, priority>)`,
183 with the usual semantics, and an `ApproxDeleteMin()` operations, ensuring:

- 184 (1) **Rank Bound.** `ApproxDeleteMin()` always returns one of the top k elements in priority
185 order.

186 (2) ***k*-Fairness.** A *priority inversion* on element $\langle \text{key}, \text{priority} \rangle$ is the event that
187 `ApproxDeleteMin()` returns a key with a *lower* priority while $\langle \text{key}, \text{priority} \rangle$ is in the
188 queue. Any element can experience at most k priority inversions before it must be returned.

189 Relaxed schedulers are quite popular in practice, as several efficient implementations have been
190 proposed and applied [36, 6, 39, 4, 18, 26, 31, 3, 35], with state-of-the-art results in the graph
191 processing domain [26, 17, 19]. A parallel line of work has attempted to provide guarantees on
192 the amount of relaxation in individual schedulers [3, 2, 34], as well as the impact of using relaxed
193 scheduling on existing iterative algorithms [1, 5]. Here, we employ the modeling of relaxed schedulers
194 used in e.g. [2, 5] for graph algorithms, but apply it to inference on graphical models.

195 3.2 Relaxed priority-based belief propagation

196 Given a Markov random field, a priority-based schedule for BP is defined by a set of *tasks*
197 T_1, T_2, \dots, T_K , each corresponding to a sequence of edge updates, and a priority function r that
198 assigns a priority $r(T_i)$ to a task based on the current state of the messages as well as possible auxiliary
199 information maintained separately. As discussed in the introduction, a non-relaxed algorithm can
200 store all tasks in a priority queue, iteratively retrieve the highest-priority task, perform all its message
201 updates, and update priorities accordingly.

202 The relaxed variant works in exactly the same way, but assuming a *k-relaxed* priority scheduler Q_k .
203 More precisely, the following steps are repeated until a fixed convergence criterion is reached:

- 204 (1) $T_i \leftarrow Q_k.\text{ApproxDeleteMin}()$ selects a task T_i among the k of highest priority in Q_k .
- 205 (2) Perform all message updates specified by the task T_i .
- 206 (3) Update the priorities for all tasks.

207 Note that tasks can be executed multiple times in this framework. In particular, we assume that the
208 priority $r(T_i)$ of a task T_i can only remain the same or increase when other tasks are executed, and
209 the only point where the priority decreases is when the task is actually executed.

210 3.3 Concurrent implementation

211 The sequential version of a priority-based schedule for belief propagation can be implemented using a
212 priority queue Q . One could map the sequential pattern directly to a parallel setting, by replacing the
213 sequential priority queue with a linearizable concurrent one. However, this may not be the best option,
214 for two reasons. First, it is challenging to build *scalable* exact priority queues [21]—the data structure
215 is inherently contended, which leads to poor cache behavior and poor performance. Second, in this
216 context, linearizability only gives the illusion of atomicity with respect to task message updates: the
217 data structure only ensures that the *task removal* is atomic, whereas the actual message updates which
218 are part of the task are not usually performed atomically together with the removal.

219 **The Multiqueue.** For this reason, in our framework, we use a *relaxed* priority scheduler, i.e. a scal-
220 able approximate priority queue called the Multiqueue [32, 3]. As the name suggests, the Multiqueue
221 is composed of m independent *exact* priority queues. To `Insert` an element, a thread picks one of
222 the exact priority queues uniformly at random, and inserts into it. To perform `ApproxDeleteMin()`,
223 the thread picks *two* of the priority queues uniformly at random, and removes the *higher priority*
224 element among their two top elements. Although very simple, this strategy has been shown to have
225 strong probabilistic rank and fairness guarantees:

226 **Theorem 1** ([3, 1]). *A Multiqueue formed of $m \geq 3$ priority queues ensures the rank and fairness*
227 *guarantees with parameter $k = O(m \log m)$, with high probability.*

228 **Our Implementation.** For our purposes, we assume that each thread i has one or a few local
229 concurrent priority queues, used to store pointers to BP-specific tasks (e.g. messages), which are
230 prioritized by an algorithm-specific function, e.g. the residual values for residual BP. We store
231 additional metadata as required by the algorithm and the graphical model. (In our experiments, we
232 use binary heaps for these priority queues, protected by locks.) To process a new task, the thread
233 selects two among all the priority queues uniformly at random, and accesses the task/message from
234 the queue whose top element has higher priority. The task is marked as *in-process* so it cannot be
235 processed concurrently by some other thread. The thread then proceeds to perform the metadata
236 updates required by the underlying variant of belief propagation, e.g., updating the message and the
237 priorities of messages from the corresponding node. The termination condition, e.g., the magnitude
238 of the largest update, is checked periodically.

239 4 Dynamics of relaxed belief propagation on trees

240 As we will see in Section 5, the relaxed priority-based belief propagation schedules yield fast
241 converge times on a wide variety of Markov random fields; specifically, the number of message
242 updates is roughly the same as for the non-relaxed version, while the running times are lower. The
243 complementary theoretical question we examine here is whether we can give analytical bounds how
244 much extra work—in terms of additional message updates—the relaxation incurs in the worst-case.

245 Unfortunately, the dynamics of even synchronous belief propagation are poorly understood on general
246 graphs, and no priority-based algorithms provide general guarantees on the convergence time. As
247 such, we can only hope to gain some limited understanding on why relaxation retains the fast
248 convergence properties of the exact priority-based schedules. Here, we present some theoretical
249 evidence suggesting that as long as a schedule does not impose long dependency chains in the
250 sequence of updates, relaxation incurs low overhead, but also that simple (non-loopy) worst-case
251 instances exist.

252 **Analytical model.** For analysis of the relaxed priority-based belief propagation, we consider the
253 formal model introduced by [5, 2] to analyze performance of iterative algorithms under relaxed
254 schedulers. Specifically, we model a relaxed scheduler Q_k as a data structure which stores pairs
255 corresponding to tasks and their priorities, with the operational semantics given in Section 3. In
256 particular, there exists a parameter k such that each `ApproxDeleteMin` returns one of the k highest
257 priority tasks in Q_k , and if a task T becomes the highest priority task in Q_k at some point during the
258 execution, then one of the next k `ApproxDeleteMin` operations returns T . (By [3, 1], our practical
259 implementation will satisfy these conditions with parameter $k = O(p \log p)$ w.h.p., where p is the
260 number of concurrent threads.) Other than satisfying these properties, we assume that the behavior of
261 Q_k can be adversarial, or randomized.

262 We model the behavior of relaxed priority-based belief propagation by investigating the number
263 of message updates needed for convergence when the algorithm is executed *sequentially* using a
264 relaxed scheduler Q_k satisfying the above constraints. This analysis reduces to a sequential game
265 between the algorithm, which queries Q_k for tasks/messages, and the scheduler, which returns
266 messages in possibly arbitrary fashion. One may think of the relaxed sequential execution as a form
267 of linearization for the actual parallel execution—reference [1] formalizes this intuition. Please see
268 the discussion at the end of this section for a practical interpretation.

269 **Relaxed BP on trees.** We now consider the behavior of relaxed residual belief propagation schedules
270 on *trees with a single source*. The setting is similar to the analysis of residual splash of Gonzalez et
271 al. [16]. Specifically, we assume that the Markov random field and the initialization of the algorithm
272 satisfies (1) The graph $G = (V, E)$ is a tree with a specified root r ; and (2) The factors of the Markov
273 random field and the initial messages are such that the residuals are zero for all messages other than
274 the outgoing messages from the root, i.e., $\text{res}(\mu_{i \rightarrow j}) = 0$ if $i \neq r$.

275 These conditions mean that residual belief propagation will start from the root, and propagate the
276 messages down the trees until propagation reaches all leaves. In particular, residual belief propagation
277 without relaxation will perform $n - 1$ message updates before convergence, updating each message
278 away from root once. While this setting is restrictive, it does model practical instances where the
279 MRF has tree-like structure, such as LDPC codes (see Section 5).

280 To characterize the dynamics on relaxed residual belief propagation on trees with a single source, we
281 observe that the algorithm can make two types of message updates:

- 282 – Updating a message with zero residual, in which case nothing happens (*a wasted update*). This
283 happens if the scheduler relaxes past the range of messages with non-zero residual.
- 284 – Updating a message $\mu_{i \rightarrow j}$ with non-zero residual, in which case the residual of $\mu_{i \rightarrow j}$ goes down
285 to zero, and the messages $\mu_{j \rightarrow k}$ for the children k of j may change their residuals to non-zero
286 values (*a useful update*).

287 It follows that each edge will get updated only once with non-zero residual. At any point of time
288 during the execution of the algorithm, we say that the *frontier* is the set of messages with non-zero
289 residual, and use $F(t)$ to denote the size of the frontier at time step t .

290 To see how the size of the frontier relates to the number message updates in relaxed residual belief
291 propagation, observe that after a useful update, we have one of the following cases:

- 292 – If $F(t) \geq k$, then the next `ApproxDeleteMin()` operation to Q_k will give an edge with non-zero
293 residual, resulting in a useful update.

294 – If $F(t) < k$, then in the worst case we need k ApproxDeleteMin() operations until we perform
295 a useful update.

296 Our main analytic result bounds the worst-case work incurred by relaxation in two concrete cases.

297 **Lemma 2.** Assume a k -relaxed scheduler Q_k for belief propagation in the tree model defined above.
298 The total number of updates performed by relaxed residual BP can be bounded as follows:

299 – **Good case: uniform expansion.** If the tree model has identical and non-deterministic edge factors
300 ψ_{ij} with $\psi_{ij}(x_i, x_j) \neq 0$ for all $\{i, j\}$, then the total number of updates performed by relaxed
301 residual BP is $n + O(Hk^2)$.

302 – **Bad case: long paths.** There exists a tree instance with height $H = o(n)$ and an adversarial
303 scheduler where relaxed residual belief propagation performs $\Omega(kn)$ message updates.

304 **Discussion.** The full analysis and illustrations are provided in Appendix A and B. To interpret the
305 results, first note that, in practice, the relaxation factor is in the order of p , the number of threads,
306 and that H is usually small (e.g., logarithmic) w.r.t. the total number of baseline updates n . Thus,
307 in the good case, the $O(k^2H)$ overhead can be seen as negligible: as p iterations occur in parallel,
308 the average time-to-completion should be $n/p + O(kH)$, which suggests almost perfect parallel
309 speedup. At the same time, our worst-case instances shows that relaxed residual BP is not a “silver
310 bullet:” there exist tree instances where it may lead to $\Omega(kn)$ message updates, i.e. asymptotically no
311 speedup. We discuss how to alter the priority function to mitigate these worst-case instances on trees
312 in the Appendix. The next section shows experimentally that such worst-case instances are unlikely.

313 5 Evaluation

314 We now empirically test the performance of the relaxed priority-based algorithms, comparing it
315 against prior work. For the experiments, we have implemented multiple priority-based algorithms
316 and instantiated them with both exact and relaxed priority schedulers.

317 **Priority-based algorithms.** We implemented several variants of sequential belief propagation,
318 among which synchronous (round-robin), residual, weight decay, and residual without lookahead.
319 These variants were briefly described in Section 2, and are specified in detail in Appendix C.
320 For residual splash, we implemented two variants. The first is the standard splash algorithm, as
321 given in [17]. The second is our own optimized version we refer to as *smart splash*, which only
322 updates messages along breadth-first-search edges during a splash operation. This variant has similar
323 convergence as the baseline residual splash algorithm, but performs fewer message updates and
324 should be more efficient. We include the following instantiations of the algorithms in the benchmarks
325 and compare them against the sequential residual algorithm. Please see Section 3.3 or Appendix C
326 for implementation details.

327 **Prior work algorithms.** We ran many possible concurrent variants for the baseline algorithms,
328 and chose the four best to we compare against our relaxed versions. For the full results, please see
329 Appendix E. First, we choose the parallel version of the standard synchronous belief propagation
330 (Synch). We omit some synchronous algorithms such as the randomized synchronous belief prop-
331 agation of Van der Merve et al. [11] since they perform consistently worse (Appendix E.2.4). We
332 also include the exact residual algorithms implemented via the coarse-grained priority queue, which
333 maintains exact priority order. Specifically, here we implemented standard the residual BP algorithm
334 (Coarse-G) and the splash algorithm of [17] (Splash) with the best value of H , which we found to
335 be 10.

336 We also include the randomized version of splash algorithm (Random Splash) proposed in the
337 journal version of the paper [16] with $H = 2$, which performed best. This algorithm uses a similar
338 idea of relaxation, but, crucially, instead of a Multiqueue scheduler, they implement a naive relaxed
339 queue where threads randomly insert and delete into p exact priority queues. While this distinction
340 may seem small, it is known [3] that this variant does not implement a k -relaxed scheduler for
341 any k , as its relaxation factor grows (diverges) as more and more operations are performed, and
342 therefore corresponds to picking tasks at random to perform. As we will see in our experimental
343 analysis, this does result in a significant difference between the number of additional (wasted) steps
344 performed relative to a relaxed priority scheduler. Finally, we note that we are the first to implement
345 this algorithm.

346 **Relaxed algorithms.** We compare the above algorithms against the algorithms we propose, i.e.
347 relaxed versions of residual belief propagation (Relaxed Residual), weight decay belief propagation
348 (Weight-Decay), residual without lookahead (Priority) and smart splash (Relaxed Smart Splash)

Input	Prior Work				Relaxed			
	Synch	Coarse-G	Splash (10)	Random Splash (2)	Residual	Weight-Decay	Priority	Smart Splash (2)
Tree	2.538x	0.265x	1.648x	2.252x	1.391x	1.282x	1.239x	2.121x
Ising	3.009x	0.801x	5.393x	11.731x	6.720x	6.276x	5.759x	14.175x
Potts	—	0.624x	1.041x	11.855x	7.454x	5.978x	5.850x	15.235x
LDPC	17.735x	1.166x	—	5.150x	13.393x	5.615x	—	10.519x

Table 1: Algorithm speedups with respect to the sequential residual algorithm. Higher is better.

Input	Prior Work				Relaxed			
	Synch	Coarse-G	Splash (10)	Random Splash (2)	Residual	Weight-Decay	Priority	Smart Splash (2)
Tree	48.000x	1.003x	16.442x	8.344x	1.020x	1.012x	3.657x	2.565x
Ising	45.006x	1.003x	9.266x	5.787x	1.058x	1.068x	1.816x	1.878
Potts	—	1.006x	9.005x	5.983x	1.068x	1.053x	1.791x	1.891x
LDPC	4.404x	1.003x	—	4.089x	1.007x	0.883x	—	0.973x

Table 2: Total updates relative to the sequential residual algorithm at 70 threads. Lower is better.

349 with the best value $H = 2$. For all these algorithms, the scheduler is a Multiqueue with 4 priority
350 queues per thread, as discussed in Section 3, which we found to work best (although other values
351 behave similarly).

352 **Methodology.** We run our experiments on four MRFs of moderate size: a binary tree of size 10^7 , an
353 Ising model [14, 20] of size $10^3 \times 10^3$, a Potts [38] of size $10^3 \times 10^3$ and the decoding of (3, 6)-LDPC
354 code [30] of size $3 \cdot 10^5$. More information about tasks and running times given in Appendix D.

355 For each pair of algorithm and model, we run each experiment five times, and average the execution
356 time and the number of performed updates on the messages. We executed on a 4-socket Intel Xeon
357 Gold 6150 2.7 GHz server with four sockets, each with 18 cores, and 512GB of RAM. The code
358 is written in Java; we use Java 11.0.5 and OpenJDK VM 11.0.5. Our code is available at <https://cutt.ly/neurips20208924>. Our code is fairly well optimized—in sequential executions it
359 outperforms the C++-based open-source framework of libDAI [25] by more than 10x, and by more
360 than 100x with multi-threading. The sizes of the inputs described in the previous paragraph are chosen
361 such that their execution is fairly long while the data still can fit into RAM.
362

363 We run the baseline algorithm on one process since it is sequential, while all other algorithms are
364 concurrent and, thus, are executed using 70 threads. The execution times (speedups) relative to
365 the sequential baseline are presented in Table 1. Each cell of the table shows how much faster the
366 corresponding algorithm works in comparison to the sequential residual one, i.e., higher is better;
367 “—” means that the execution did not converge within a reasonable limit of time.

368 **Results.** See Table 1 for the speedups versus the baseline, on 70 threads. (For ablation studies,
369 see Appendix E.) On trees, the fastest algorithm is, predictably, the synchronous one, since on
370 tree-like models with small diameter D it performs only approximately $O(D)$ times more updates in
371 comparison to the sequential baseline, while being almost perfectly parallelizable. So, it works well
372 on perfect binary tree as our Tree model, but works much worse on the chain graphs. On Ising and
373 Potts models, the best algorithm is Relaxed Smart Splash (RSS) with $H = 2$. The algorithm closest to
374 it is Random Splash with $H = 2$, which is 20 – 30% slower. For LDPC decoding, which is a tree-like
375 model, the best-performing is again the Synchronous algorithm. We note the good performance of the
376 relaxed residual algorithm, as well as of RSS, and the relatively poor performance of Random Splash,
377 due to high numbers of wasted updates. Examining Table 2, we notice in general the relatively low
378 number of wasted updates for relaxed algorithms. In summary, the choice of algorithm can depend on
379 the model; however, one may choose Relaxed Smart Splash since it performs well on all our models.

380 6 Discussion

381 We have investigated the use of relaxed schedulers in the context of the classic belief propagation
382 algorithm for inference on graphical model, and have shown that this approach leads to an efficient
383 family of algorithms, which improve upon the previous state-of-the-art non-relaxed parallelization
384 approaches in our experiments. Overall, our relaxed implementations, either Relaxed Residual or
385 Relaxed Smart Splash, have state-of-the-art performance in multithreaded regimes, making them a
386 good generic choice for any belief propagation task.

387 For future work, we highlight two possible directions. First is to extend our theoretical analysis
388 to cover more types of instances; however, as we have seen, the structure of belief propagation
389 schedules can be quite complicated, and the challenge is the figure out a proper framework for more
390 general analysis. Second possible direction is extending our empirical study to a massively-parallel,
391 multi-machine setting.

392 **Broader impact**

393 As this work is focused on speeding up existing inference techniques and does not focus on a specific
394 application, the main benefit is enabling belief propagation applications to process data sets more
395 efficiently, or enable use of larger data sets. We do not expect direct negative societal consequences
396 to follow from our work, though we note that as with all heuristic machine learning techniques, there
397 is an inherent risk of misinterpreting the results or ignoring biases in the data if proper care is not
398 taken in application of the methods. However, such risks exist regardless of our work.

399 **References**

- 400 [1] Dan Alistarh, Trevor Brown, Justin Kopinsky, Jerry Z. Li, and Giorgi Nadiradze. Distributionally
401 linearizable data structures. In *Proceedings of the 30th on Symposium on Parallelism in*
402 *Algorithms and Architectures*, SPAA '18, pages 133–142, New York, NY, USA, 2018. ACM.
- 403 [2] Dan Alistarh, Trevor Brown, Justin Kopinsky, and Giorgi Nadiradze. Relaxed schedulers can
404 efficiently parallelize iterative algorithms. In *Proceedings of the 2018 ACM Symposium on*
405 *Principles of Distributed Computing*, PODC '18, pages 377–386, New York, NY, USA, 2018.
406 ACM.
- 407 [3] Dan Alistarh, Justin Kopinsky, Jerry Li, and Giorgi Nadiradze. The power of choice in priority
408 scheduling. In *Proceedings of the ACM Symposium on Principles of Distributed Computing*,
409 PODC '17, pages 283–292, New York, NY, USA, 2017. ACM.
- 410 [4] Dan Alistarh, Justin Kopinsky, Jerry Li, and Nir Shavit. The SprayList: A scalable relaxed
411 priority queue. In *20th ACM SIGPLAN Symposium on Principles and Practice of Parallel*
412 *Programming*, PPOPP 2015, San Francisco, CA, USA, 2015. ACM.
- 413 [5] Dan Alistarh, Giorgi Nadiradze, and Nikita Koval. Efficiency guarantees for parallel incremental
414 algorithms under relaxed schedulers. In *The 31st ACM Symposium on Parallelism in Algorithms*
415 *and Architectures*, SPAA '19, page 145–154, New York, NY, USA, 2019. Association for
416 Computing Machinery.
- 417 [6] Dmitry Basin, Rui Fan, Idit Keidar, Ofer Kiselov, and Dmitri Perelman. CAFÉ: Scalable
418 task pools with adjustable fairness and contention. In *Proceedings of the 25th International*
419 *Conference on Distributed Computing*, DISC'11, pages 475–488, Berlin, Heidelberg, 2011.
420 Springer-Verlag.
- 421 [7] Guy E Blelloch, Yan Gu, Julian Shun, and Yihan Sun. Parallelism in randomized incremental
422 algorithms. In *Proceedings of the 28th ACM Symposium on Parallelism in Algorithms and*
423 *Architectures*, pages 467–478. ACM, 2016.
- 424 [8] Andres I Vila Casado, Miguel Griot, and Richard D Wesel. Informed dynamic scheduling
425 for belief-propagation decoding of LDPC codes. In *2007 IEEE International Conference on*
426 *Communications*, pages 932–937. IEEE, 2007.
- 427 [9] Gregory F. Cooper. The computational complexity of probabilistic inference using Bayesian
428 belief networks. *Artificial Intelligence*, 42(2):393–405, 1990.
- 429 [10] Paul Dagum and Michael Luby. Approximating probabilistic inference in Bayesian belief
430 networks is NP-hard. *Artificial Intelligence*, 60(1):141–153, 1993.
- 431 [11] Mark Van der Merwe, Vinu Joseph, and Ganesh Gopalakrishnan. Message scheduling for
432 performant, many-core belief propagation. In *Proceedings of the IEEE High Performance*
433 *Extreme Computing Conference (HPEC 2019)*, 2019.
- 434 [12] Laxman Dhulipala, Guy Blelloch, and Julian Shun. Julienne: A framework for parallel graph
435 algorithms using work-efficient bucketing. In *Proceedings of the 29th ACM Symposium on*
436 *Parallelism in Algorithms and Architectures*, SPAA '17, pages 293–304, New York, NY, USA,
437 2017. ACM.
- 438 [13] Laxman Dhulipala, Guy E. Blelloch, and Julian Shun. Theoretically efficient parallel graph
439 algorithms can be fast and scalable. In *30th on Symposium on Parallelism in Algorithms and*
440 *Architectures (SPAA 2018)*, pages 393–404, 2018.

- 441 [14] Gal Elidan, Ian McGraw, and Daphne Koller. Residual belief propagation: informed scheduling
442 for asynchronous message passing. In *Proceedings of the 22nd Conference on Uncertainty in*
443 *Artificial Intelligence (UAI 2006)*, pages 165–173, 2006.
- 444 [15] Robert G. Gallager. Low-density parity-check codes. *IRE Transactions on Information Theory*,
445 8(1):21–28, 1962.
- 446 [16] Joseph Gonzalez, Yucheng Low, and Carlos Guestrin. Residual splash for optimally parallelizing
447 belief propagation. In *Proceedings of the 12th International Conference on Artificial Intelligence*
448 *and Statistics (UAI 2009)*, volume 5, pages 177–184, 2009.
- 449 [17] Joseph E Gonzalez, Yucheng Low, Haijie Gu, Danny Bickson, and Carlos Guestrin. Powergraph:
450 distributed graph-parallel computation on natural graphs. In *10th USENIX Symposium on*
451 *Operating Systems Design and Implementation (OSDI '12)*, pages 17–30, 2012.
- 452 [18] Andreas Haas, Michael Lippautz, Thomas A. Henzinger, Hannes Payer, Ana Sokolova,
453 Christoph M. Kirsch, and Ali Sezgin. Distributed queues in shared memory: multicore per-
454 formance and scalability through quantitative relaxation. In *Computing Frontiers Conference,*
455 *CF'13, Ischia, Italy, May 14 - 16, 2013*, pages 17:1–17:9, 2013.
- 456 [19] Mark C Jeffrey, Suvinay Subramanian, Cong Yan, Joel Emer, and Daniel Sanchez. Unlocking
457 ordered parallelism with the swarm architecture. *IEEE Micro*, 36(3):105–117, 2016.
- 458 [20] Christian Knoll, Michael Rath, Sebastian Tschiatschek, and Franz Pernkopf. Message schedul-
459 ing methods for belief propagation. In *Machine Learning and Knowledge Discovery in*
460 *Databases (ECML PKDD 2015)*, pages 295–310, 2015.
- 461 [21] Andrew Lenharth, Donald Nguyen, and Keshav Pingali. Priority queues are not good concurrent
462 priority schedulers. In *European Conference on Parallel Processing*, pages 209–221. Springer,
463 2015.
- 464 [22] Mu Li, David G Andersen, Jun Woo Park, Alexander J Smola, Amr Ahmed, Vanja Josifovski,
465 James Long, Eugene J Shekita, and Bor-Yiing Su. Scaling distributed machine learning with the
466 parameter server. In *11th USENIX conference on Operating Systems Design and Implementation*
467 *(OSDI'14)*, page 583–598, 2014.
- 468 [23] Ji Liu and Stephen J Wright. Asynchronous stochastic coordinate descent: Parallelism and
469 convergence properties. *SIAM Journal on Optimization*, 25(1):351–376, 2015.
- 470 [24] Yucheng Low, Joseph E. Gonzalez, Aapo Kyrölä, Danny Bickson, Carlos E. Guestrin, and
471 Joseph Hellerstein. GraphLab: A new framework for parallel machine learning. In *26th*
472 *Conference on Uncertainty in Artificial Intelligence (UAI 2010)*, page 340–349, 2010.
- 473 [25] Joris M. Mooij. libDAI: A free and open source C++ library for discrete approximate inference
474 in graphical models. *Journal of Machine Learning Research*, 11:2169–2173, August 2010.
- 475 [26] Donald Nguyen, Andrew Lenharth, and Keshav Pingali. A lightweight infrastructure for
476 graph analytics. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems*
477 *Principles, SOSP '13*, pages 456–471, New York, NY, USA, 2013. ACM.
- 478 [27] Judea Pearl. Reverend Bayes on inference engines: A distributed hierarchical approach. In
479 *Proceedings of the Second AAAI Conference on Artificial Intelligence (AAAI 1982)*, page
480 133–136. AAAI Press, 1982.
- 481 [28] Benjamin Recht, Christopher Re, Stephen Wright, and Feng Niu. Hogwild: A lock-free
482 approach to parallelizing stochastic gradient descent. In *Advances in neural information*
483 *processing systems*, pages 693–701, 2011.
- 484 [29] Thomas J. Richardson and Rüdiger L. Urbanke. The capacity of low-density parity-check codes
485 under message-passing decoding. *IEEE Transactions on information theory*, 47(2):599–618,
486 2001.
- 487 [30] Thomas J. Richardson and Rüdiger L. Urbanke. *Modern coding theory*. Cambridge university
488 press, 2008.

- 489 [31] Hamza Rihani, Peter Sanders, and Roman Dementiev. Brief announcement: MultiQueues:
490 Simple relaxed concurrent priority queues. In *Proceedings of the 27th ACM Symposium on*
491 *Parallelism in Algorithms and Architectures*, SPAA '15, pages 80–82, New York, NY, USA,
492 2015. ACM.
- 493 [32] Hamza Rihani, Peter Sanders, and Roman Dementiev. Brief announcement: Multiqueues:
494 Simple relaxed concurrent priority queues. In *Proceedings of the 27th ACM symposium on*
495 *Parallelism in Algorithms and Architectures*, pages 80–82, 2015.
- 496 [33] Dan Roth. On the hardness of approximate reasoning. *Artificial Intelligence*, 82(1):273–302,
497 1996.
- 498 [34] Adones Rukundo, Aras Atalar, and Philippas Tsigas. Monotonically Relaxing Concurrent
499 Data-Structure Semantics for Increasing Performance: An Efficient 2D Design Framework. In
500 *33rd International Symposium on Distributed Computing (DISC 2019)*, volume 146 of *Leibniz*
501 *International Proceedings in Informatics (LIPIcs)*, pages 31:1–31:15, Dagstuhl, Germany, 2019.
502 Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
- 503 [35] Konstantinos Sagonas and Kjell Winblad. A contention adapting approach to concurrent ordered
504 sets. *Journal of Parallel and Distributed Computing*, 2017.
- 505 [36] Nir Shavit and Itay Lotan. Skiplist-based concurrent priority queues. In *Parallel and Distributed*
506 *Processing Symposium, 2000. IPDPS 2000. Proceedings. 14th International*, pages 263–268.
507 IEEE, 2000.
- 508 [37] Julian Shun, Guy E. Blelloch, Jeremy T. Fineman, and Phillip B. Gibbons. Reducing contention
509 through priority updates. In *Proceedings of the Twenty-fifth Annual ACM Symposium on*
510 *Parallelism in Algorithms and Architectures*, SPAA '13, pages 152–163, New York, NY, USA,
511 2013. ACM.
- 512 [38] Charles Sutton and Andrew McCallum. Improved dynamic schedules for belief propagation. In
513 *Proceedings of the 23rd Conference on Uncertainty in Artificial Intelligence (UAI 2007)*, pages
514 376–383, 2007.
- 515 [39] Martin Wimmer, Jakob Gruber, Jesper Larsson Träff, and Philippas Tsigas. The lock-free
516 k -LSM relaxed priority queue. In *20th ACM SIGPLAN Symposium on Principles and Practice*
517 *of Parallel Programming (PPoPP 2015)*, pages 277–278, 2015.
- 518 [40] Jonathan S. Yedidia, William T. Freeman, and Yair Weiss. Understanding belief propagation
519 and its generalizations. In *Exploring Artificial Intelligence in the New Millenium, chapter 8*,
520 pages 239—269. Morgan Kaufmann, 2003.

521 **A Analysis**

522 **Lemma 2.** Assume a k -relaxed scheduler Q_k for belief propagation in the tree model defined above.
 523 The total number of updates performed by relaxed residual BP can be bounded as follows:

- 524 – **Good case: uniform expansion.** If the tree model has identical and non-deterministic edge factors
- 525 ψ_{ij} with $\psi_{ij}(x_i, x_j) \neq 0$ for all $\{i, j\}$, then the total number of updates performed by relaxed
- 526 residual BP is $n + O(Hk^2)$.
- 527 – **Bad case: long paths.** There exists a tree instance with height $H = o(n)$ and an adversarial
- 528 scheduler where relaxed residual belief propagation performs $\Omega(kn)$ message updates.

529 *Proof. Good case: uniform expansion.* As the first case, we consider the tree model in the case
 530 where the edge factors ψ_{ij} are identical for all edges and not *deterministic*, i.e. $\psi_{ij}(x_i, x_j) \neq 0$ for
 531 all $\{i, j\}$. Let us say that the *level* of a message $\mu_{i \rightarrow j}$ is ℓ if the distance from i to the root r is ℓ . The
 532 conditions we imposed our Markov random field, together with the update rule (2), imply that the
 533 residuals of the messages are decreasing in the level ℓ of the message, and all messages on level ℓ will
 534 have the same residual when they are in the frontier. This means that residual schedule will prefer
 535 updating messages on lower levels first.

536 Now consider the progress of the relaxed residual belief propagation on this tree; let H denote the
 537 height of the tree. Now assume that all messages on levels $0, 1, \dots, \ell - 1$ have had a useful update,
 538 and consider how many wasted updates we can make in the worst case before all messages on level ℓ
 539 have been processed. Let f denote the number of message on level ℓ still in the frontier:

- 540 – While $f \geq k$, there are at least k messages of level ℓ on the frontier. Since they have the highest
- 541 residual out of the messages in the frontier, each update is a useful update of a message on level ℓ .
- 542 – When $f < k$, there can be updates that do not hit messages on level ℓ , which can possibly be
- 543 wasted updates. However, the highest-priority messages are still from level ℓ , so every k th update
- 544 will hit a message on level ℓ by the guarantees of the scheduler. Thus, in $(k - 1)f = O(k^2)$
- 545 updates, all remaining messages on level ℓ have been processed.

546 Since there can be at most $n - 1$ useful updates, and the number of levels is $H - 1$, the total number
 547 of updates performed by relaxed residual belief propagation is $n + O(Hk^2)$.

548 **Bad case: long paths.** A simple example where relaxed residual belief propagation performs poorly
 549 is a path. That is, if our underlying tree is a path of length n with a root at one end, then relaxed
 550 residual belief propagation can perform $\Omega(kn)$ message updates in the worst case. However, the path
 551 has height $H = n$, so one might ask if there is a general upper bound of form $n + O(Hk^2)$ on trees
 552 without restricting the edge factors as in our previous example.

553 Unfortunately, turns out that without the restrictions above, we can construct examples of trees with
 554 height $H = o(n)$ where relaxed residual belief propagation still performs $\Omega(kn)$ message updates
 555 (see Figure 2 for an illustration):

- 556 (1) Start with a path of length \sqrt{n} , with a root at one end.
- 557 (2) Attach a new path of length \sqrt{n} to each vertex.
- 558 (3) For each remaining degree-2 node in the graph, attach a single new node to it.

This construction results in a 3-regular rooted tree with $\Theta(n)$ nodes and depth $H = O(\sqrt{n})$.
 Finally, we choose the edge factors so that residuals on the side paths are larger than the residuals
 on the main path, so residual belief propagation will prefer following the side paths first.

One can now observe that under the adversarial model for the relaxed scheduler, the adversary
 can select the execution of the relaxed scheduler so that the frontier size never exceeds 4. That
 is, adversary forces the algorithm to process the graph one side path at a time, wasting $k - 1$ steps
 between each useful update.

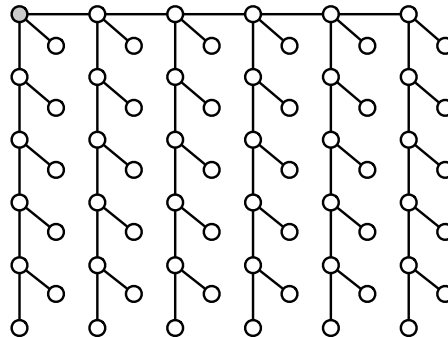


Figure 2: Example of the tree where relaxed residual belief propagation performs poorly.

Finally, we note that the same construction can be generalized to obtain instances with similar relaxation overhead and diameter $O(n^{1/c})$ for larger constants $c < k$, by simply working with paths of length $n^{1/c}$ and repeating the path attachment step c times. □

559 *Remark 3.* As suggested by the above examples, one might consider changing the priority function
 560 to preferentially select messages closer to the source. This can lead to improved work guarantees for
 561 the relaxed schedule. Indeed, we discuss one concrete example in Section B, where we show how
 562 to relax the optimal schedule on trees. However, it is not straightforward to construct such priority
 563 functions so that they also make sense on general graphs, which can have non-monotonic potentials
 564 and cycles.

565 B Optimal schedule on trees

On trees, the belief propagation gives exact marginals under any schedule that updates each edge infinitely often. However, there is an optimal schedule that updates each message exactly once, requiring $O(n)$ message updates [27]. Assume the tree has a fixed root v :

- 566 (1) In the first phase, all messages towards the root are updated starting from the leaves; each
 567 message is updated only after all its predecessors have been updated.
- 568 (2) In the second phase, all messages away from the root are update starting from the root.

569 This schedule can be modeled in the priority-based scheduling framework as follows:

- 570 (1) Initially, the outgoing messages at leaf nodes have priority n , and all other messages have
 571 priority 0.
- 572 (2) When message is updated with non-zero priority, its priority is changed to 0.
- 573 (3) Once all messages $\mu_{k \rightarrow i}$ for $k \in N(i) \setminus \{j\}$ have been updated once with non-zero priority,
 574 the message $\mu_{i \rightarrow j}$ changes to priority to minimum of update priorities of the incoming edges
 575 minus one.

576 This priority function can clearly be implemented by keeping a constant amount of extra information
 577 per message. When the above schedule is executed with an exact scheduler, the algorithm will update
 578 each message once with non-zero priority before considering any messages with zero priority, and by
 579 following the analysis of [27], one can see that the algorithm has converged at that point.

580 Similarly, in the relaxed version of the schedule, the algorithm has converged once all messages have
 581 been updated once with non-zero priority. In addition, some messages may be updated multiple times
 582 with priority 0; we call these *wasted* updates, and the updates done while the message has non-zero
 583 priority *useful* updates.

584 **Claim 4.** *The relaxed version of the optimal schedule on trees performs $O(n + k^2 H)$ message*
 585 *updates, where H is the height of the tree.*

586 *Proof.* For the purposes of analysis, assign messages into buckets B_1, B_2, \dots, B_n so that bucket B_ℓ
 587 contains the messages that will have their useful update done with priority ℓ . One can observe that
 588 the update priority of message $\mu_{i \rightarrow j}$ is the $n - d$, where d is the maximum distance from node i to a
 589 leaf using a path that does not cross edge $\{i, j\}$. Since this is bounded by the diameter of the tree,
 590 there are at most $2H$ non-empty buckets.

591 Assume that all messages in buckets $B_n, B_{n-1}, \dots, B_{\ell+1}$ have been already had a useful update. We
 592 now show that in there can be at most k^2 wasted updates before all messages in B_ℓ have had a useful
 593 update. Since all earlier buckets have been processed, all messages in B_ℓ have either already had a
 594 useful update, or have priority ℓ . Let b be the number of messages remaining in bucket B_ℓ :

- 595 – While $b \geq k$, there are at least k messages with priority ℓ , so each update is a useful update
 596 of a message in B_ℓ
- 597 – When $b < k$, there can be wasted updates. However, since buckets $B_n, B_{n-1}, \dots, B_{\ell+1}$
 598 have had all useful updates, the top elements in the schedule will be from bucket B_ℓ , and

599 thus by the guarantees of the scheduler, there can be at most $k - 1$ wasted updates before
600 the top element is processed. Thus, in $b(k - 1) = O(k^2)$ updates, all remaining messages
601 of B_ℓ will have their useful update.

602 By an inductive argument, all non-empty buckets have been processed after $O(k^2H)$ wasted updates,
603 so the total number of updates is $O(n + k^2H)$. \square

604 C Algorithms

605 C.1 Asynchronous belief propagation

606 Starting with Elidan et al. [14], there has been a line of research arguing that *asynchronous* or *iterative*
607 schedules for belief propagation tend to converge more reliably and with fewer message updates than
608 the synchronous schedule. In particular, the practical work has focused on developing schedules that
609 attempt to iteratively perform ‘the most useful’ update at each step; the most prominent of these
610 algorithms is the *residual belief propagation* of Elidan et al. [14], with other proposals aiming to
611 address the shortcomings of residual belief propagation in various cases.

612 **Residual belief propagation.** Given a current state of messages, let $\mu'_{i \rightarrow j}$ denote the message we
613 would obtain by applying the message update rule (2) to message $\mu_{i \rightarrow j}$. In residual belief propagation,
614 the priority of a message is given by the *residual* $\text{res}(\mu_{i \rightarrow j})$ of a message $\mu_{i \rightarrow j}$, defined as

$$\text{res}(\mu_{i \rightarrow j}) = \|\mu'_{i \rightarrow j} - \mu_{i \rightarrow j}\|, \quad (4)$$

615 where $\|\cdot\|$ is an arbitrary norm; in this work, we assume L^2 norm is used unless otherwise specified.
616 That is, the residual of a message corresponds to amount of change that would happen if message $\mu_{i \rightarrow j}$
617 would be updated. Note that this means that residual belief propagation performs *lookahead*, that is,
618 the algorithm precomputes the future updates before applying them to the state of the algorithm.

619 **Weight decay belief propagation.** *Weight decay belief propagation* of [20] is a variant of residual
620 belief propagation that penalizes message priorities for repeated updates. That is, let $m(\mu_{i \rightarrow j})$ denote
621 how many times message $\mu_{i \rightarrow j}$ has been updated by the algorithm, and let $\text{res}(\mu_{i \rightarrow j})$ denote the
622 residual of a message as above. The priority function of weight decay belief propagation is

$$r(\mu_{i \rightarrow j}) = \frac{\text{res}(\mu_{i \rightarrow j})}{m(\mu_{i \rightarrow j})}.$$

623 The motivation behind this weight decay scheme is that empirical observations suggest that one
624 possible failure mode of residual belief propagation is getting stuck in cycles with large residuals; the
625 weight decay prioritizes other edges in cases where this happens.

626 **Residual without lookahead.** Another variant of residual belief propagation is the lookahead-
627 avoiding belief propagation of [38]. As the name implies, this algorithm does not perform the exact
628 residual computation using (4), but instead approximates the residuals indirectly, with the aim of
629 reducing the computational cost of priority updates.

630 Informally, the basic idea is that for each message $\mu_{i \rightarrow j}$, we track the amount other incoming
631 messages at node i have changed since the last update of $\mu_{i \rightarrow j}$, and use this to define the priority of
632 updating $\mu_{i \rightarrow j}$. The actual approximation in the algorithm uses a slightly different notion of residual
633 from (4), so we refer to [38] for full details.

634 C.2 Parallel belief propagation

635 As discussed above, the question of parallelizing belief propagation is fairly poorly understood.
636 The synchronous schedule is trivially parallelizable by performing updates within each round in
637 parallel, but the improved converge properties of the iterative schedules cannot easily be translated to
638 parallel setting. There have been recent proposals that aim to bridge this gap in an ad-hoc manner by
639 designing custom algorithms for specific parallel computation settings.

640 **Residual splash.** The *residual splash* belief propagation [16] is a vertex-based algorithm inspired
641 by residual belief propagation. The residual splash algorithm was initially designed for MapReduce
642 computation, and it aims to have larger individual tasks while retaining a similar structure to residual
643 belief propagation.

644 Specifically, the residual splash algorithm works by defining a priority function over nodes of the
 645 Markov random field, and selecting the next node to process in a strict priority order. For the selected
 646 node, the algorithm performs a *splash* operation that propagates information within distance H in the
 647 graph; in practice, this results in threads performing larger individual tasks at once, offsetting the cost
 648 of accessing the strict scheduler.

649 In detail, the priority of for nodes is given by the *node residual*, defined as

$$\text{res}(i) = \max_{j \in N(i)} \text{res}(\mu_{j \rightarrow i}).$$

650 Given a *depth parameter* H , the splash operation at node i is defined by following sequence of
 651 message updates:

- 652 (1) Construct a BFS tree T of depth H rooted at node i .
- 653 (2) In the reverse BFS order on T —starting from leaves—process all nodes in T , updating all
 654 outgoing messages for each node processed.
- 655 (3) Repeat the previous step in BFS order, i.e., starting from the root.

656 In other words, this process gather all available information at radius H from the selected node, and
 657 propagates it to all nodes within the radius.

658 **Randomized synchronous belief propagation.** Van der Merve et al. [11] proposed a parallelization
 659 scheme for belief propagation on GPUs, mixing the structure of synchronous and residual belief
 660 propagation. Their algorithm considers all messages at once in global rounds, and performs the
 661 following filter-and-select steps before computing the message updates:

- 662 (1) Filter out all messages whose residuals are below the convergence threshold.
- 663 (2) Out of the remaining messages, select a p fraction of messages uniformly at random to
 664 update.

665 Alternatively, the process can perform the algorithm on per-node basis, using node residuals as in the
 666 residual splash algorithm.

667 The fraction p is adjusted on the fly based on the convergence of the algorithm, preferring a low
 668 value if the algorithm is converging slowly, and a high value if it is converging fast. Concretely, the
 669 selection scheme for p used by [11] is to set $p = 1$ if the number of messages above the convergence
 670 threshold decreased by at least 10% in the last round, and set it to a smaller fixed value otherwise.

671 We note that the randomized synchronous algorithm is particularly well suited for GPU use, as
 672 the filter-and select steps can be efficiently implemented on GPUs. However, as shown by our
 673 experimental study, this strategy is not efficient on a subset of real-world models, when ported to
 674 CPU. Conversely, as discussed by the authors of [11], the dynamic priority-based strategy we propose
 675 would be hard to implement efficiently on GPUs, due to its irregular structure.

676 D Models

677 We run our experiments on four Markov random fields models.

678 **Trees.** As a simple base case, we consider a simple tree model similar to the analytical setting in
 679 Section 4. The underlying graph is a full binary tree on 10 millions vertices, and the other parameters
 680 are set up as follows:

- 681 – All variables are binary, i.e. the domain is $\{0, 1\}$ for each variable.
- 682 – Vertex factors are $(0.1, 0.9)$ for the root and $(0.5, 0.5)$ for all other vertices.
- 683 – Edge factors are $\psi_{ij}(x, y) = \begin{cases} 1, & x = y \\ 0, & x \neq y \end{cases}$ for all edges.

684 As discussed in Section 4, these choices create a setup where the belief propagation has to propagate
 685 information from the root to all other nodes. Thus, under an optimal schedule, the total number of
 686 performed updates is be equal to $10^7 - 1$. Since we know that all algorithms will converge on this
 687 model, we run the algorithms until exact convergence.

688 **Ising and Potts models.** Ising and Potts models are Markov random fields defined over an $n \times n$
689 grid graph, arising from applications in statistical physics. Both of Ising [14, 20] and Potts [38]
690 models were used in prior work as test case, and in general they offer a class of good test instances,
691 as they both exhibit complex cyclic propagations and are easy to generate.

692 For the parameters of the models, we mostly follow prior work in the setup. As the underlying graph,
693 we use a $10^3 \times 10^3$ grid graph to get instances where the effects of parallelization are clearly visible.
694 For the Ising model, we select the factors similarly to [14, 20]:

- 695 – The variable domain is $\{-1, 1\}$ for all variables.
- 696 – Vertex factors are $\psi_i(x) = e^{\beta_i x}$.
- 697 – Edge factors are $\psi_{ij}(x, y) = e^{\alpha_{ij} xy}$.
- 698 – The parameters α_{ij} and β_i are chosen uniformly at random from $[-1, 1]$.

699 For the Potts model, we select the factors following [38]:

- 700 – The variable domain is $\{0, 1\}$ for all variables.
- 701 – Vertex factors are $\psi_i(x) = \begin{cases} e^{\beta_i}, & x = 1 \\ 1, & x = 0 \end{cases}$.
- 702 – Edge factors are $\psi_{ij}(x, y) = \begin{cases} e^{\alpha_{ij}}, & x = y \\ 1, & x \neq y \end{cases}$.
- 703 – The parameters α_{ij} and β_i are chosen uniformly at random from $[-2.5, 2.5]$.

704 For both Ising and Potts models, we set the convergence threshold to 10^{-5} . That is, we terminate
705 algorithm once all task have priority below this threshold.

706 **LDPC codes.** Finally, we generate Markov random fields corresponding to the $(3, 6)$ -LDPC (*low*
707 *density parity check code* [15]) decoding. LDPC decoding is one of the more successful application
708 of belief propagation. We consider a simple version of LDPC decoding task where convergence
709 guarantees exist [29]. However, we stress that coding theory is its own extensive research area, and
710 far more optimized codes and decoding algorithms exist in practice—we simply use LDPC decoding
711 to observe the comparative scaling behavior of our implementations on instances where synchronous
712 belief propagation is guaranteed to converge. For a more detailed background on LDPC decoding
713 and other aspects of coding theory, refer e.g. to the book [30].

714 More precisely, we consider $(3, 6)$ -LDPC decoding over a binary symmetric channels. Informally,
715 a $(3, 6)$ -LDPC code is a $(3, 6)$ -regular bipartite graph, where each degree 3 node corresponds to a
716 binary *variable* and each degree 6 node corresponds to a *constraint* of form $x_{i_1} + x_{i_2} + \dots + x_{i_6} = 0$
717 over the neighboring variables $x_{i_1}, x_{i_2}, \dots, x_{i_6}$. Each sequence of variables that satisfies the all the
718 constraints is *codeword* of the code. The basic setup is then that we send a codeword over a *channel*
719 that flips each bit with probability ε , and the receiver will run belief propagation and use results of
720 marginalization to infer the original codeword.

721 For our experiments, we build a $(3, 6)$ -LDPC instance with 300 000 variable nodes and 150 000
722 constraint nodes by selecting a random $(3, 6)$ -regular bipartite graph, and initialize the node factors
723 corresponding to the all-zero codeword sent over binary symmetric channel with error probability
724 $\varepsilon = 0.07$. Under these conditions, belief propagation is guaranteed to correctly decode the instance
725 with high probability [29]; indeed, all the algorithms that converged decoded the codeword correctly
726 in our experiments. The codeword length was again selected to get roughly comparable baseline
727 running times as for the other instances.

728 Concretely, we get Markov random field where the underlying graph is a random bipartite graph with
729 450 000 nodes. For each variable node i , let $x_i \in \{0, 1\}$ be the ‘transmitted’ value of the variable,
730 randomly generated to be 1 with probability ε and 0 otherwise. The factors have the following
731 structure:

- 732 – The domains of variable nodes are binary domains $\{0, 1\}$. For the constraint nodes, the
733 domain is $\{0, 1\}^6$ —different bit masks of length 6.
- 734 – The node factors for variable nodes are

$$\psi_i(y) = \begin{cases} 1 - \varepsilon, & y = x_i \\ \varepsilon, & y \neq x_i. \end{cases}$$

Input	Residual	Prior Work						Relaxed				
		Synch	CG	S 2	S 10	RS 2	RS 10	Residual	WD	Priority	RSS 2	RSS 10
Tree	1.30 min	2.538x	0.265x	0.608x	1.648x	2.252x	2.241x	1.391x	1.282x	1.239x	2.121x	2.110x
Ising	2.76 min	3.009x	0.801x	0.609x	5.393x	11.731x	13.512x	6.720x	6.276x	5.759x	14.175x	10.337x
Potts	3.02 min	—	0.624x	0.484x	1.041x	11.855x	12.854x	7.454x	5.978x	5.850x	15.235x	11.091x
LDPC	4.62 min	17.735x	1.166x	—	—	5.150x	—	13.393x	5.615x	—	10.519x	—

Table 3: Algorithm speedups with respect to the sequential residual algorithm. Higher is better.

Input	Residual	Prior Work						Relaxed				
		Synch	CG	S 2	S 10	RS 2	RS 10	Residual	WD	Priority	RSS 2	RSS 10
Tree	10M	48.000x	1.003x	8.658x	16.442x	8.344x	15.197x	1.020x	1.012x	3.657x	2.565x	5.027x
Ising	25.3M	45.006x	1.003x	5.719x	9.266x	5.787x	10.232x	1.058x	1.068x	1.816x	1.878x	6.147x
Potts	30M	—	1.006x	5.903x	9.005x	5.983x	9.109x	1.068x	1.053x	1.791x	1.891x	6.328x
LDPC	7.23M	4.404x	1.003x	—	—	4.089x	—	1.007x	0.883x	—	0.973x	—

Table 4: Total updates relative to the sequential residual algorithm at 70 threads. Lower is better.

735 For the constraint nodes, the node factor $\psi_c(y)$ is equal to the number of ones in $y \in \{0, 1\}^6$
736 modulo 2; this effectively penalizes any value that does not satisfy the constraint.
737 – Edge factors $\psi_{ic}(x, y)$ is one if the corresponding bit in the $y \in \{0, 1\}^6$ equals $x \in \{0, 1\}$,
738 and is zero otherwise.

739 For the LDPC instances, we set the convergence threshold to 10^{-2} to ensure fast convergence; this
740 approximates the behavior of actual LDPC decoders.

741 E Experiments

742 In this section, we provide an additional study on the evaluation of the algorithms. At first, we give
743 the extended results of running the algorithms on the moderate size inputs chosen in the main body of
744 the paper. Unfortunately, due to the reasonably high running time it was impossible to make enough
745 points for the plots to reason about the general effect of the parallelization. Thus, we execute the
746 algorithms on a little bit smaller inputs.

747 E.1 Moderate size inputs

748 To present all the executed algorithms in the table, we shrink the abbreviations a little bit: Coarse-
749 Grained now becomes CG, Splash becomes S, Random Splash becomes RS, Relaxed Residual
750 rests Residual, Weight-Decay becomes WD, Relaxed Priority becomes Priority, and, finally,
751 Relaxed Smart Splash becomes RSS. Table 3 contains the execution times (speedups) of the
752 algorithms relative to the sequential baseline. Table 4 contains the number of updates performed by
753 the algorithms in compare to the number of updates performed by the sequential baseline. The results
754 do not differ much from the ones presented in the main body of the paper. The only notable thing is
755 that Random Splash with $H = 10$ it performs better on Ising and Potts model than Random Splash
756 with $H = 2$. However, we chose Random Splash with $H = 2$ as the best one, since Random Splash
757 with $H = 10$ does not finish on LDPC input. Nevertheless, Relaxed Smart Splash outperforms
758 Random Splash with both settings of H .

759 E.2 Small size inputs

760 In this subsection, we decrease the size of the inputs. Now, Tree model maintains a tree of size 10^6 ,
761 Ising and Potts models are built on top of 300×300 grid graph, and, finally, LDPC model is set
762 up with 30 000 length of the input vector. In general, we simply reduce the sizes of the models by
763 approximately 10.

764 E.2.1 Scaling

765 **How to read the plots.** There are two types of plots per each model: the first shows the execution
766 time of the algorithms, while the other one shows the number of updates performed. On the x axis
767 we have the number of threads the algorithms were run on, while on the y axis we have: the time in
768 seconds (for time plots) and the number of updates (for update plots). The dashed lines on the plots
769 correspond to the algorithms that use a relaxed scheduler, while the others use either no concurrent
770 scheduler, or an exact priority queue.

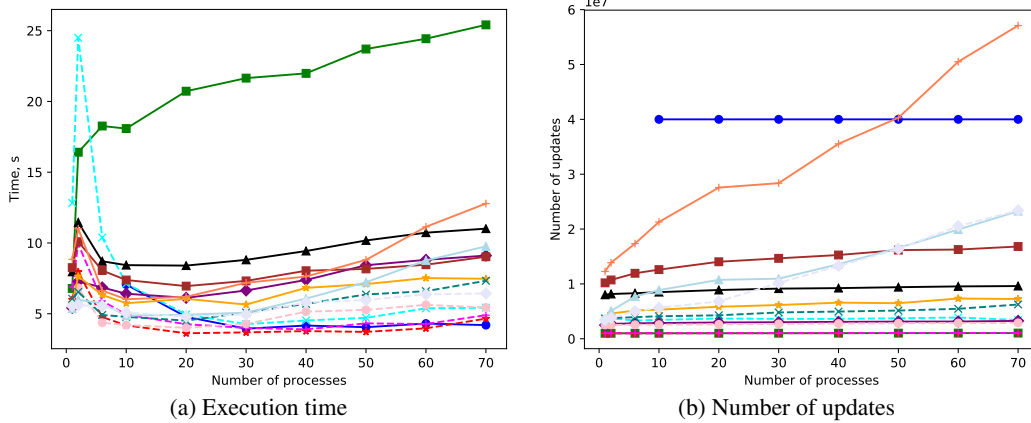
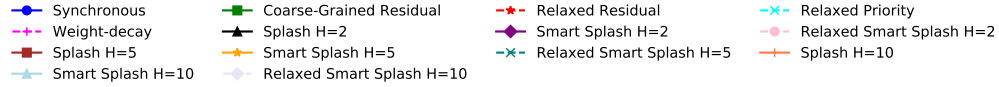


Figure 3: The results of the evaluation of the algorithms on the Tree model

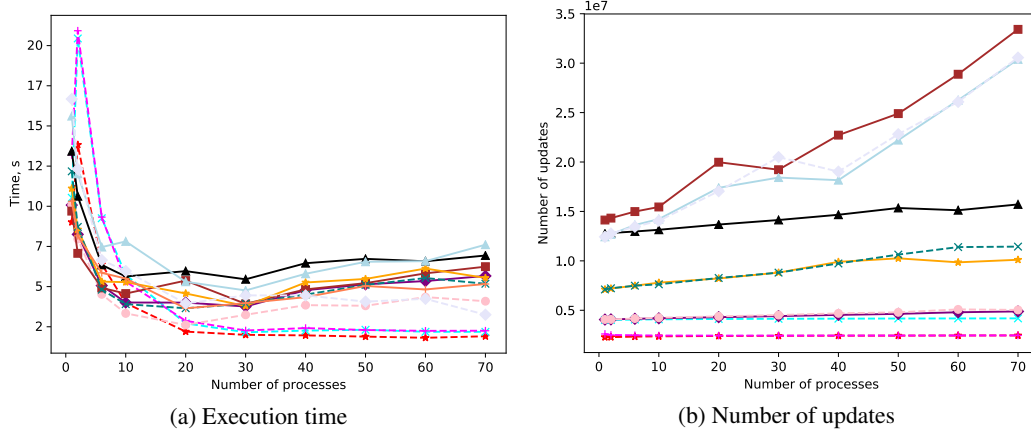


Figure 4: The results of the evaluation of the algorithms on Ising model

771 Whenever we have omitted algorithms from the plots or display incomplete data, this indicates poor
 772 performance for that algorithm on the metric displayed on the graph: either the algorithm did not
 773 converge or the values exceed the limit of the plot.

774 **Tree model.** As one can observe on the time plot (Figure 3a), the three algorithms with the best
 775 scaling on the tree instance are the synchronous belief propagation, relaxed residual and the weight-
 776 decay algorithm. For the relaxed algorithms, this mirrors our theoretical analysis from Section 4: as
 777 can be seen from Figure 3b, the relaxation incurs very low overhead in terms of additional updates,
 778 while the overhead from parallelization is also low. By contrast, the exact residual belief propagation
 779 performs exactly the minimum number of updates needed, but scales very badly due to the contention
 780 on the priority queue.

781 We note that on the tree instance, the synchronous belief propagation also scales very well when
 782 parallelized. The amount of work can be split evenly between the threads, and only $O(\log n)$
 783 synchronous rounds are required for convergence.

784 **Ising and Potts model.** Ising and Potts models represent more challenging instances with lots of
 785 cycles, and are generally thought to be more representative of hard general graph instances for belief
 786 propagation. As can be seen in Figures 4a and 5a, relaxed algorithms perform consistently well on
 787 these instances, with relaxed residual belief propagation giving consistently the fastest convergence.
 788 These are followed by the exact splash algorithms, which generally perform slightly worse; however,

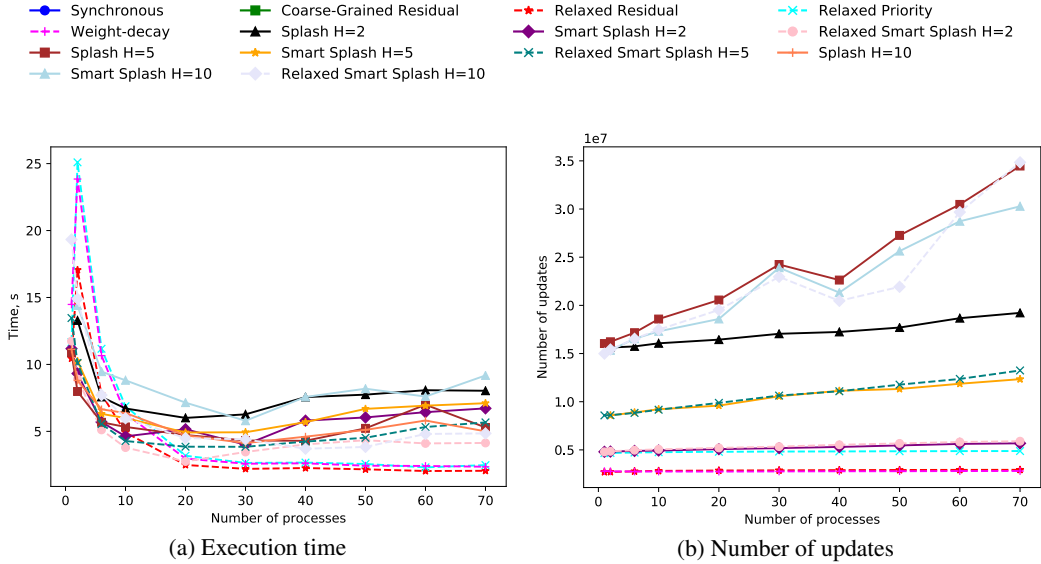


Figure 5: The results of the evaluation of the algorithms on Potts model

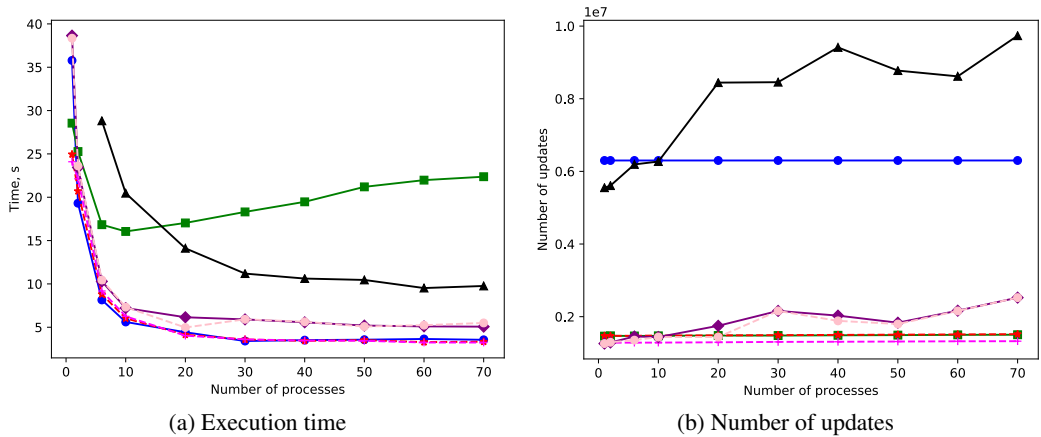


Figure 6: The results of the evaluation of the algorithms on decoding LDPC code

789 the scaling seems to be somewhat sensitive to the choice of the parameter H . Both the synchronous
 790 and exact residual belief propagation are omitted, as the former did not consistently converge, and
 791 the latter was very slow.

792 An interesting insight is that the exact variants of splash and smart splash do not converge at all in
 793 single-threaded executions for some values of the parameter H , but always converge on two and more
 794 threads. Similarly, synchronous belief propagation, which has a fixed schedule, does not converge. By
 795 contrast, relaxed smart splash converged under all parameter values. We conjecture that this is due to
 796 the phenomenon observed by [20]: exact priority-based algorithms may get stuck in non-convergent
 797 cyclic schedules, and injecting randomness into the schedule may help the algorithm to ‘escape’ these
 798 situations. In particular, relaxation to the priority queue, i.e., sometimes executing low-priority items,
 799 can provide a such source of randomness. Similarly, an increase in the number of threads leads to the
 800 relaxation of the algorithm even for exact schedulers, as several messages are processed in parallel,
 801 not only the best one. Thus, we empirically observe that the randomness in the relaxation might help
 802 belief propagation to avoid bad cyclic schedules and, therefore, converge.

803 **LDPC model.** There are five algorithms that perform similarly (Figure 6a): synchronous belief
 804 propagation, relaxed residual belief propagation, the weight decay algorithm, relaxed smart splash
 805 with $H = 2$ and, finally, smart splash with $H = 2$. The other algorithms did not converge within our
 806 five minutes time limit per experiment.

		Message updates			
	Threads	Tree	Ising	Potts	LDPC
Exact	1	1000000	2279000	2700000	1464000
Relaxed	1	+0.14%	+0.11%	-0.01%	+0.55%
	2	+0.26%	+0.24%	+0.37%	+0.57%
	6	+0.56%	+2.50%	+2.70%	+0.64%
	10	+0.92%	+3.71%	+4.45%	+1.05%
	20	+2.08%	+5.27%	+5.87%	+1.41%
	30	+2.90%	+6.10%	+6.56%	+1.87%
	40	+3.48%	+6.52%	+7.39%	+2.35%
	50	+5.04%	+6.83%	+7.92%	+2.83%
	60	+4.96%	+7.39%	+8.28%	+3.20%
70	+5.74%	+7.71%	+8.53%	+3.70%	

Table 5: Number of additional message updates performed by relaxed residual belief propagation compared to exact residual belief propagation.

807 We note that synchronous belief propagation performs very well on this instance. This is not surprising,
808 as standard belief propagation is known to perform well in LDPC decoding. Generally speaking,
809 the necessary propagation chains seem to be very short on LDPC instances, and the synchronous
810 algorithm parallelizes well in such cases.

811 E.2.2 The effects of relaxation

812 In Table 5, we measure how many more updates the relaxed residual algorithm needs to perform
813 in comparison to the number of updates performed by the standard sequential residual algorithm,
814 denoted as “baseline”. We count the total number of updates only approximately: we check the
815 convergence condition only after every 1000 iterations.

816 The left column indicates whether it is a baseline algorithm or the number of threads for relaxed
817 residual belief propagation. The other columns present the numbers for each model we consider.
818 Each cell contains the corresponding number of updates and how many more updates the relaxed
819 version of the algorithm executed (percentage).

820 On one process, relaxed residual performs more updates than the baseline does, except in the case
821 of the Potts model. It is expected since our algorithm uses relaxed Multiqueue instead of the strict
822 priority queue. Moreover, as expected the overhead on the number of updates in comparison to the
823 baseline increases with the number of threads. This is again due to the relaxation of the priority
824 queue—recall that we allocate $4\times$ more queues than threads. Interestingly, this overhead is limited
825 even on 70 threads—its maximum value is 9% maximum. This explains the good performance of our
826 algorithm: we reduce the contention by relaxing accesses to the priority queue, while at the same
827 time the total number of updates does not increase significantly.

828 E.2.3 Relaxed versus Non-Relaxed Algorithms

829 In Table 6, we analyze the speedups obtained by the relaxed residual algorithm relative to the best-
830 performing non-relaxed alternative across models and thread counts. We notice that our algorithm
831 outperforms the alternatives in most of the cases, often by a large margin—the highest speedup is of
832 $2.85\times$, whereas the highest slow-down is of $0.47x$. Both occur on the Potts model, which is generally
833 the most difficult instance in our tests. Overall, the combination of our relaxed scheduling framework
834 combined with the standard residual belief propagation is clearly the algorithm of choice at high
835 thread counts, where it consistently outperforms the alternatives; on the other hand, relaxed residual
836 also performs reasonably well on a single thread, making it a consistently good choice all across the
837 board.

838 E.2.4 Random Synchronous Algorithm

839 In Table 7, we present the execution time of random synchronous algorithm on 70 threads (Random
840 Synch 70) with different values of $lowP = 0.1, 0.4$ and 0.7 , where the parameter $lowP$ controls the

Threads	Speedup			
	Tree	Ising	Potts	LDPC
1	0.89x	1.08x	1.04x	1.14x
2	0.75x	0.51x	0.47x	1.13x
6	1.20x	0.77x	0.73x	1.17x
10	1.16x	1.01x	0.94x	1.20x
20	1.36x	1.66x	1.89x	1.49x
30	1.38x	1.88x	1.82x	1.65x
40	1.61x	2.21x	1.90x	1.62x
50	1.91x	2.67x	2.36x	1.48x
60	1.89x	2.66x	2.85x	1.55x
70	1.61x	2.71x	2.44x	1.52x

Table 6: Speedup of relaxed residual belief propagation versus the best non-relaxed alternative on different thread counts. We note that overhead of parallelization can overcome the benefits on small thread counts, as seen in the scaling experiments.

841 random selection fraction p in steps where the algorithm is converging slowly (see Section C.2). We
842 compare it with the execution time of two baselines: Synchronous algorithm on 70 threads (Synch
843 70) and Relaxed Residual on one process (RR 1). Cells with ‘—’ indicate executions that either take
844 more than five minutes to run or simply do not converge.

845 To summarize, we did not include the execution time of random synchronous algorithm in the scaling
846 plots since it exceeds the execution time of one of the baselines in all cases.

Algorithm		Running time (s)			
		Tree	Ising	Potts	LDPC
Synch 70		4.088	—	—	3.504
RR 1		5.579	9.012	10.583	25.663
Random Synch 70	$lowP = 0.1$	37.052	62.629	—	28.543
	$lowP = 0.4$	8.420	20.396	—	7.269
	$lowP = 0.7$	6.306	12.581	—	4.791

Table 7: Randomized synchronous algorithm versus baselines.