

# Parallel Non-blocking Deterministic Algorithm for Online Topic Modeling

Oleksandr Frei<sup>1</sup> and Murat Apishev<sup>2</sup>

<sup>1</sup> Moscow Institute of Physics and Technology, [oleksandr.frei@gmail.com](mailto:oleksandr.frei@gmail.com)

<sup>2</sup> National Research University Higher School of Economics, [great-mel@yandex.ru](mailto:great-mel@yandex.ru)

**Abstract.** In this paper we present a new asynchronous algorithm for learning additively regularized topic models and discuss the main architectural details of our implementation. The key property of the new algorithm is that it behaves in a fully deterministic fashion, which is typically hard to achieve in a non-blocking parallel implementation. The algorithm had been recently implemented in the BigARTM library (<http://bigartm.org>). Our new algorithm is compatible with all features previously introduced in BigARTM library, including multimodality, regularizers and scores calculation. While the existing BigARTM implementation compares favorably with the alternative packages such as Vowpal Wabbit or Gensim, the new algorithm brings yet further improvements in CPU utilization, memory usage, and spends even less time to achieve the same perplexity.

**Keywords:** probabilistic topic modeling, Probabilistic Latent Semantic Analysis, Latent Dirichlet Allocation, Additive Regularization of Topic Models, stochastic matrix factorization, EM-algorithm, online learning, asynchronous and parallel computing, BigARTM.

## 1 Introduction

Topic models [1] is a powerful machine learning technology for statistical text analysis that has been widely used in text mining, information retrieval, network analysis and other areas [2]. Today a lot of research efforts around topic models is devoted to distributed implementations of *Latent Dirichlet Allocation* (LDA) [4], a specific Bayesian topic model that uses Dirichlet conjugate prior. This lead to numerous implementations such as AD-LDA [7], PLDA [8] and PLDA+ [9], all designed to run on a big cluster. Largest topic models of web scale can reach millions of topics and vocabulary words, yielding Big Data models with trillions of parameters [3]. Yet not all researchers and application are dealing with so large web-scale collections, and they require an efficient implementation that can run on a powerful workstation or even laptop. Such implementations are very useful, as shown by popular open-source packages Vowpal Wabbit [11], Gensim [10] and Mallet [12], which are neither distributed nor sometimes even multi-threaded.

Scaling down a distributed algorithm can be challenging. LightLDA [3] is a major step in this direction, however it focuses only on the LDA model. Our goal is to develop a flexible framework that can learn a wide variety of topic models.

BigARTM [17] is an open-source library for regularized multimodal topic modeling of large collections. BigARTM is based on a novel technique of additive regularized topic models (ARTM) [14,15,16,18], which gives flexible multi-criteria approach to probabilistic topic modeling. ARTM includes all popular models such as LDA [4], PLSA [5], and many others. Key feature of ARTM is that it provides a cohesive framework that allows users to combine different topic models that previously did not fit together.

BigARTM is proven to be very fast comparing to the alternative packages. According to [17], BigARTM runs approx. 10 times faster comparing to Gensim [10] and twice faster than Vowpal Wabbit (VW) [11] in a single thread. With multiple threads BigARTM wins even more as it scales linearly up to at least 16 threads. In this paper we address the remaining limitations of the library, including few performance bottlenecks and fixing non-deterministic behavior of the Online algorithm.

The rest of the paper is organized as follows. In section 2 we introduce basic notation. In sections 3 and 4 we summarize offline and online algorithms for learning ARTM models. In sections 5 and 6 we discuss two asynchronous non-blocking modifications of the online algorithm. In section 7 we compare the internal architecture of BigARTM library between versions 0.6 and 0.7. In section 8 we report results of our experiments on large datasets. In section 9 we discuss advantages, limitations and open problems of BigARTM.

## 2 Notation

Let  $D$  denote a finite set (collection) of texts and  $W$  denote a finite set (vocabulary) of all terms from these texts. Let  $n_{dw}$  denote the number of occurrences of a term  $w \in W$  in a document  $d \in D$ ;  $n_{dw}$  values form a sparse matrix of size  $|W| \times |D|$ , known as *bag-of-words* representation of the collection.

Given an  $(n_{dw})$  matrix, a probabilistic topic model finds two matrices:  $\Phi = \{\phi_{wt}\}$  and  $\Theta = \{\theta_{td}\}$ , of sizes  $|W| \times |T|$  and  $|T| \times |D|$  respectively, where  $|T|$  is a user-defined number of *topics* in the model. Matrices  $\Phi$  and  $\Theta$  provide a compressed representation of the  $(n_{dw})$  matrix:

$$n_{dw} \approx n_d \sum_{t \in T} \phi_{wt} \theta_{td}, \text{ for all } d \in D, w \in W,$$

where  $n_d = \sum_{w \in W} n_{dw}$  denotes the total number of terms in a document  $d$ .

To learn  $\Phi$  and  $\Theta$  from  $(n_{dw})$  an additively-regularized topic model (ARTM) maximizes the log-likelihood, regularized via an additional penalty term  $R(\Phi, \Theta)$ :

$$\sum_{d \in D} \sum_{w \in W} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}. \quad (1)$$

Regularization penalty  $R(\Phi, \Theta)$  may incorporate external knowledge of the expert about the collection. With no regularization ( $R = 0$ ) it corresponds to PLSA [5]. Many Bayesian topic models, including LDA [4], can be represented as special cases of ARTM with different regularizers  $R$ , as shown in [15,16].

---

**Algorithm 1:** ProcessDocument( $d, \Phi$ )

---

**Input:** document  $d \in D$ , matrix  $\Phi = (\phi_{wt})$ ;  
**Output:** matrix  $(\tilde{n}_{wt})$ , vector  $\theta_{td}$ ;

- 1 initialize  $\theta_{td} := \frac{1}{|T|}$  for all  $t \in T$ ;
- 2 **repeat**
- 3      $p_{tdw} := \mathop{\text{norm}}_{t \in T}(\phi_{wt}\theta_{td})$  for all  $w \in d$  and  $t \in T$ ;
- 4      $\theta_{td} := \mathop{\text{norm}}_{t \in T}(\sum_{w \in d} n_{dw}p_{tdw} + \theta_{td} \frac{\partial R}{\partial \theta_{td}})$  for all  $t \in T$ ;
- 5 **until**  $\theta_d$  converges;
- 6  $\tilde{n}_{wt} := n_{dw}p_{tdw}$  for all  $w \in d$  and  $t \in T$ ;

---

In [14] it is shown that the local maximum  $(\Phi, \Theta)$  of the problem (1) satisfies the following system of equations:

$$p_{tdw} = \mathop{\text{norm}}_{t \in T}(\phi_{wt}\theta_{td}); \quad (2)$$

$$\phi_{wt} = \mathop{\text{norm}}_{w \in W} \left( n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right); \quad n_{wt} = \sum_{d \in D} n_{dw}p_{tdw}; \quad (3)$$

$$\theta_{td} = \mathop{\text{norm}}_{t \in T} \left( n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right); \quad n_{td} = \sum_{w \in d} n_{dw}p_{tdw}; \quad (4)$$

where operator  $\mathop{\text{norm}}_{i \in I} x_i = \frac{\max\{x_i, 0\}}{\sum_{j \in I} \max\{x_j, 0\}}$  transforms a vector  $(x_i)_{i \in I}$  to a discrete distribution,  $n_{wt}$  counters represent term frequency of word  $w$  in topic  $t$ .

Learning of  $\Phi$  and  $\Theta$  from (2)–(4) can be done by EM-algorithm, which starts from a random values in  $\Phi$  and  $\Theta$ , and iterates E-step (2) and M-steps (3),(4) until convergence. In the sequel we discuss several variations of such EM-algorithm, which are all based on the above formulas but differ in the way how operations are ordered and grouped together.

In addition to plain text many collections has other metadata, such as authors, class or category labels, date-time stamps, or even associated images, audio or video clips, usage data, etc. In [18] this data can be represented as *modalities*, where the overall vocabulary  $W$  is split into  $M$  subsets  $W = W^1 \sqcup \dots \sqcup W^M$ , one subset per modality, and in (3) matrix  $\Phi$  is normalized independently within each modality. Incorporating modalities into topic model improves its quality and makes it applicable for classification, cross-modal retrieval, or making recommendations. In the sequel we list all algorithms for one modality, but our implementation in BigARTM supports the general case.

### 3 Offline algorithm

Offline ARTM (Alg. 2) relies on subroutine ProcessDocument (Alg. 1), which corresponds to equations (2) and (4) from the solution of the ARTM optimization problem (1). ProcessDocument requires a fixed  $\Phi$  matrix and a vector  $n_{dw}$  of

---

**Algorithm 2: Offline ARTM**

---

**Input:** collection  $D$ ;  
**Output:** matrix  $\Phi = (\phi_{wt})$ ;

- 1 initialize  $(\phi_{wt})$ ;
- 2 create batches  $D := D_1 \sqcup D_2 \sqcup \dots \sqcup D_B$ ;
- 3 **repeat**
- 4      $(n_{wt}) := \sum_{b=1, \dots, B} \sum_{d \in D_b} \text{ProcessDocument}(d, \Phi)$ ;
- 5      $(\phi_{wt}) := \text{norm}_{w \in W}(n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}})$ ;
- 6 **until**  $(\phi_{wt})$  converges;

---

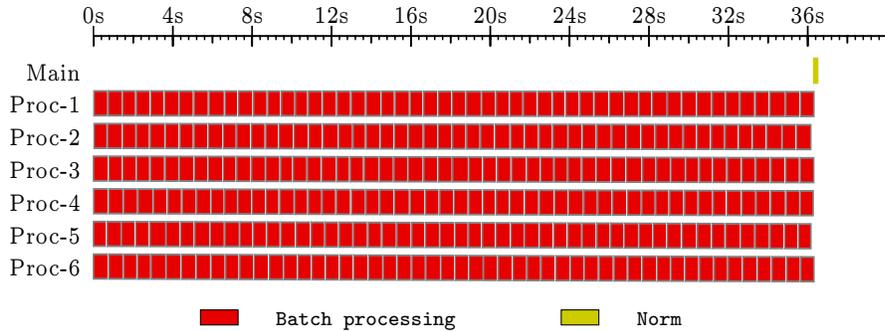


Fig. 1. Gantt chart for Offline ARTM (Alg. 2)

term frequencies for a given document  $d \in D$ , and as a result it returns a topical distribution  $\theta_{td}$  for the document, and a matrix  $(\hat{n}_{wt})$  of size  $|d| \times |T|$ , where  $|d|$  gives the number of distinct terms in the document. The `ProcessDocument` might be also useful as a separate routine which finds  $\theta_{td}$  distribution for a new document, but in the Offline algorithm it is rather used as a building block in an iterative EM-algorithm that learns the  $\Phi$  matrix.

Offline algorithm performs scans over the collection, calling `ProcessDocument` for each document  $d \in D$  from the collection, and then aggregating the resulting  $(\hat{n}_{wt})$  matrices into the final  $(n_{wt})$  matrix of size  $|W| \times |T|$ . After each scan it recalculates  $\Phi$  matrix according to the equation (3).

At step 2 we split collection  $D$  into batches  $(D_b)$ . This step is not strictly necessary for Offline algorithm, and it rather reflects an internal implementation detail. For performance reasons the outer loop over batches  $b = 1, \dots, B$  is parallelized across multiple threads, and within each batch the inner loop over documents  $d \in D_b$  is executed in a single thread. Each batch is stored in a separate disk file on disk to allow out-of-core streaming of the collection. For typical collections it is reasonable to have around 1000 documents per batch, however for ultimate performance we encourage users to experiment with this parameter. Too small batches can cause disk IO overhead due to lots of small reads, while

too large batches will result in bigger tasks that will not be distributed evenly across computation threads.

Note that  $\theta_{td}$  values appear only within `ProcessDocument` subroutine. This leads to efficient memory usage because the implementation never stores the entire theta matrix at any given time. Instead,  $\theta_{td}$  values are recalculated from scratch on every pass through the collection.

Fig. 1 shows a Gantt chart of the Offline algorithm. Here and in the sequel Gantt charts are built for a single EM-iteration on NYTimes dataset<sup>3</sup> ( $|D| = 300\text{K}$ ,  $|W| = 102\text{K}$ ) with  $|T| = 16$  topics. `ProcessBatch` boxes correspond to the time spent in processing an individual batch. The final box `Norm`, executed on the main thread, correspond to the time spent in the step 4 in Alg. 2 where  $n_{wt}$  counters are normalized to produce a new  $\Phi$  matrix.

## 4 Online algorithm

The Online ARTM (Alg. 3) generalizes the Online variational Bayes algorithm, suggested in [6] for the LDA model. Online ARTM improves the convergence rate of the Offline ARTM by re-calculating matrix  $\Phi$  after every  $\eta$  batches. To simplify the notation we introduce a trivial subroutine

$$\text{ProcessBatches}(\{D_b\}, \Phi) = \sum_{D_b} \sum_{d \in D_b} \text{ProcessDocument}(d, \Phi)$$

that aggregates the output of `ProcessDocument` across a given set of batches at a constant  $\Phi$  matrix. Here the split of the collection  $D := D_1 \sqcup D_2 \sqcup \dots \sqcup D_B$  into batches plays far more significant role than in the Offline algorithm, because different splitting algorithmically affects the result. At step 6 the new  $n_{wt}^{i+1}$  values are calculated as a convex combination of the old values  $n_{wt}^i$  and the value  $\hat{n}_{wt}^i$  produced on the recent batches. Old counters  $n_{wt}^i$  are discounted by a factor  $(1 - \rho_i)$ , which depends on the iteration number. A common strategy is to use  $\rho_i = (\tau_0 + i)^{-\kappa}$ , where typical values for  $\tau_0$  are between 64 and 1024, for  $\kappa$  — between 0.5 and 0.7.

As in the Offline algorithm, the outer loop over batches  $D_{\eta(i-1)+1}, \dots, D_{\eta i}$  is executed concurrently across multiple threads. The problem with this approach is that all threads have no useful work to do during steps 5-7 of the Online algorithm. The threads can not start processing the next batches because a new version of  $\Phi$  matrix is not ready yet. As a result the CPU utilization stays low, and the run-time Gantt chart of the Online algorithm typically looks like in Fig. 2. Boxes `Even batch` and `Odd batch` both correspond to step 4, and indicate the version of the  $\Phi^i$  matrix (even  $i$  or odd  $i$ ). `Merge` correspond to the time spend in merging  $n_{wt}$  with  $\hat{n}_{wt}$ . `Norm` is, as before, the time spent to normalize  $n_{wt}$  counters into the new  $\Phi$  matrix.

In the next two sections we present asynchronous modifications of the online algorithm that results in better CPU utilization. The first of them (`Async ARTM`)

<sup>3</sup> <https://archive.ics.uci.edu/ml/datasets/Bag+of+Words>

---

**Algorithm 3: Online ARTM**


---

**Input:** collection  $D$ , parameters  $\eta, \tau_0, \kappa$ ;  
**Output:** matrix  $\Phi = (\phi_{wt})$ ;  
1 create batches  $D := D_1 \sqcup D_2 \sqcup \dots \sqcup D_B$ ;  
2 initialize  $(\phi_{wt}^0)$ ;  
3 **for all** update  $i = 1, \dots, \lfloor B/\eta \rfloor$   
4      $(\hat{n}_{wt}^i) := \text{ProcessBatches}(\{D_{\eta(i-1)+1}, \dots, D_{\eta i}\}, \Phi^{i-1})$ ;  
5      $\rho_i := (\tau_0 + i)^{-\kappa}$ ;  
6      $(n_{wt}^i) := (1 - \rho_i) \cdot (n_{wt}^{i-1}) + \rho_i \cdot (\hat{n}_{wt}^i)$ ;  
7      $(\phi_{wt}^i) := \text{norm}_{w \in W}(n_{wt}^i + \phi_{wt}^{i-1} \frac{\partial R}{\partial \phi_{wt}})$ ;

---

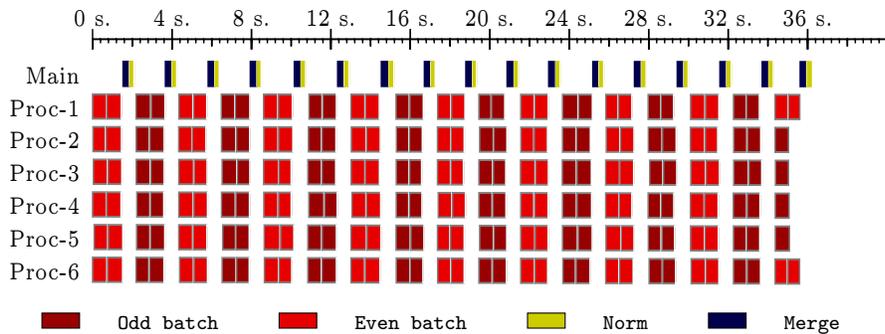


Fig. 2. Gantt chart for Online ARTM (Alg. 3)

has non-deterministic behavior and few performance bottlenecks. The second algorithm (DetAsync ARTM) addresses these problems.

## 5 Async: asynchronous online algorithm

The Async algorithm was implemented in BigARTM v0.6 as described in [17]. Basic idea is to trigger asynchronous execution of the Offline algorithm and store the resulting  $\hat{n}_{wt}$  matrices into a queue. Then, whenever the number of elements in the queue becomes a multiple of  $\eta$ , the Async algorithm performs steps 5-7 of the Online ARTM (Alg. 3). For performance reasons merging of the  $\hat{n}_{wt}$  counters happens in a background by a dedicated *Merger* thread.

First problem of the Async algorithm is that it does not define the order in which  $\hat{n}_{wt}$  are merged. This order is usually different from the original order of the batches, and typically it changes from run to run. This affects the final  $\Phi$  matrix which also changes from run to run.

Another issue with Async algorithm is that queuing  $\hat{n}_{wt}$  counters may considerably increase the memory usage, and also lead to performance bottlenecks in the *Merger* thread. In some cases the execution of the Async algorithm is

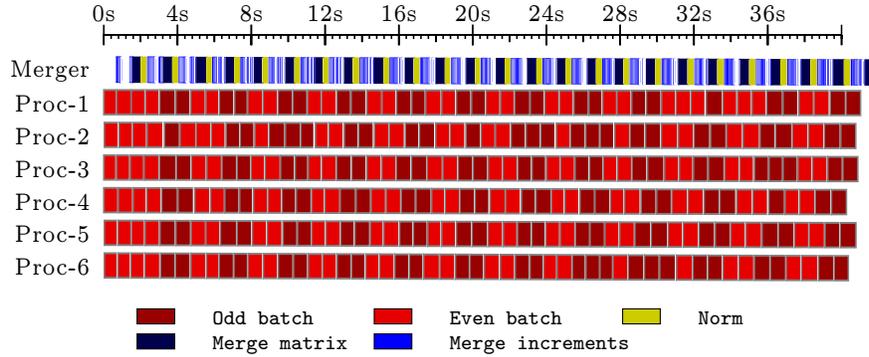


Fig. 3. Gantt chart for Async ARTM from BigARTM v0.6 — normal execution

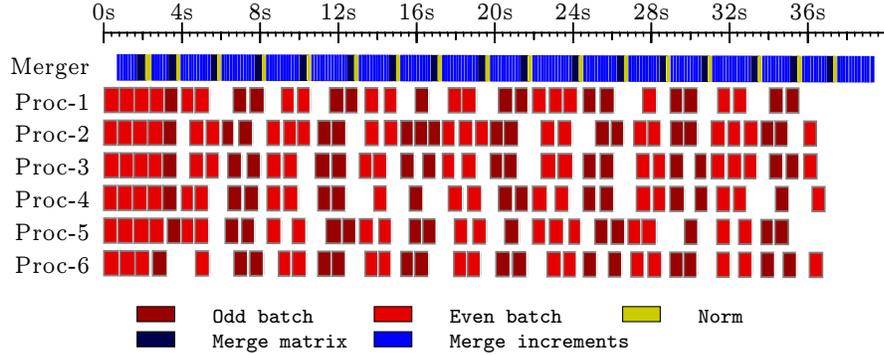


Fig. 4. Gantt chart for Async ARTM from BigARTM v0.6 — performance issues

as efficient as for the Offline algorithm, as shown on Fig. 3. However, certain combination of the parameters (particularly, small batch size or small number of iterations in `ProcessDocument`'s inner loop 2-5) might overload the merger thread. Then the Gantt chart may look as on Fig. 4, where most threads are waiting because there is no space left in the queue to place  $n_{wt}$  counters.

In the next section we resolve the aforementioned problems by introducing a new `DetAsync` algorithm, which has an entirely deterministic behavior and achieves high CPU utilization without requiring user to tweak the parameters.

## 6 DetAsync: deterministic asynchronous online algorithm

`DetAsync ARTM` (Alg. 4) is based on two new routines, `AsyncProcessBatches` and `Await`. The former is equivalent to `ProcessBatches`, except that it just queues the task for an asynchronous execution and returns immediately. Its output is a future object (for example, an `std::future` from C++11 standard), which can be

---

**Algorithm 4: DetAsync ARTM**


---

**Input:** collection  $D$ , parameters  $\eta, \tau_0, \kappa$ ;  
**Output:** matrix  $\Phi = (\phi_{wt})$ ;  
1 create batches  $D := D_1 \sqcup D_2 \sqcup \dots \sqcup D_B$ ;  
2 initialize  $(\phi_{wt}^0)$ ;  
3  $F^1 := \text{AsyncProcessBatches}(\{D_1, \dots, D_\eta\}, \Phi^0)$ ;  
4 **for all** update  $i = 1, \dots, \lfloor B/\eta \rfloor$   
5     **if**  $i \neq \lfloor B/\eta \rfloor$  **then**  
6          $F^{i+1} := \text{AsyncProcessBatches}(\{D_{\eta i+1}, \dots, D_{\eta(i+\eta)}\}, \Phi^{i-1})$ ;  
7          $(\hat{n}_{wt}^i) := \text{Await}(F^i)$ ;  
8          $\rho_i := (\tau_0 + i)^{-\kappa}$ ;  
9          $(n_{wt}^i) := (1 - \rho_i) \cdot (n_{wt}^{i-1}) + \rho_i \cdot (\hat{n}_{wt}^i)$ ;  
10          $(\phi_{wt}^i) := \text{norm}_{w \in W}(n_{wt}^i + \phi_{wt}^{i-1} \frac{\partial R}{\partial \phi_{wt}})$ ;

---

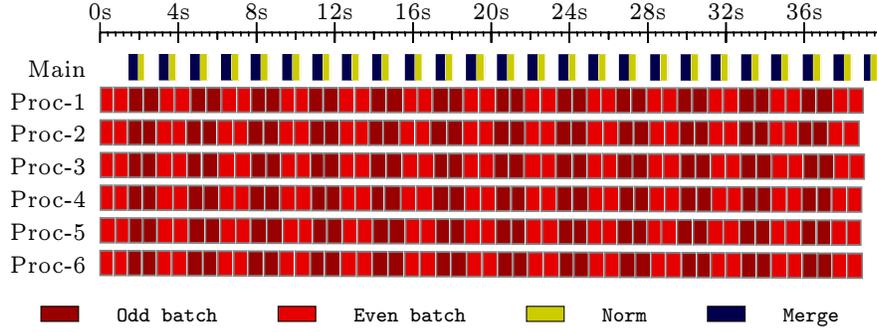


Fig. 5. Gantt chart for DetAsync ARTM (Alg. 4)

later passed to `Await` in order to get the actual result, e.g. in our case the  $\hat{n}_{wt}$  values. In between calls to `AsyncProcessBatches` and `Await` the algorithm can perform some other useful work, while the background threads are calculating the  $(\hat{n}_{wt})$  matrix.

To calculate  $\hat{n}_{wt}^{i+1}$  it uses  $\Phi^{i-1}$  matrix, which is one generation older than  $\Phi^i$  matrix used by the Online algorithm. This adds an extra “offset” between the moment when  $\Phi$  matrix is calculated and the moment when it is used, and as a result gives the algorithm additional flexibility to distribute more payload to computation threads. Steps 3 and 5 of the algorithm are just technical tricks to implement the “offset” idea.

Adding an “offset” should negatively impact the convergence of the DetAsync algorithm comparing to the Online algorithm. For example, in `AsyncProcessBatches` the initial matrix  $\Phi^0$  is used twice, and the two last matrices  $\Phi^{\lfloor B/\eta \rfloor - 1}$  and  $\Phi^{\lfloor B/\eta \rfloor}$  will not be used at all. On the other hand the asynchronous algorithm gives

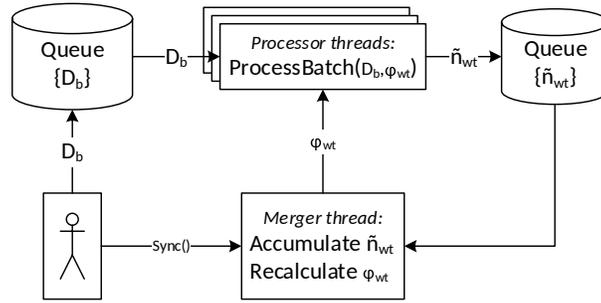


Fig. 6. Diagram of BigARTM components (old architecture)

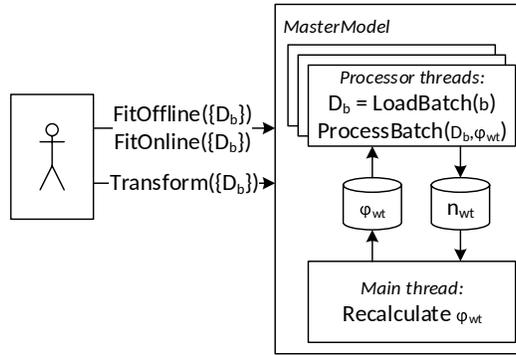


Fig. 7. Diagram of BigARTM components (new architecture)

better CPU utilization, as clearly shown by the Gantt chart from Fig. 5. This tradeoff between convergence and CPU utilization will be evaluated in section 8.

## 7 Implementation

The challenging part for the implementation is to aggregate the  $\hat{n}_{wt}$  matrices across multiple batches, given that they are processed in different threads. The way BigARTM solves this challenge was changed between versions v0.6 (Fig. 6) and v0.7 (Fig. 7).

In the old architecture the  $\hat{n}_{wt}$  matrices were stored in a queue, and then aggregated by a dedicated *Merger thread*. In the new architecture we removed Merger thread, and  $\hat{n}_{wt}$  are written directly into the final  $n_{wt}$  matrix concurrently from all processor threads. To synchronize the write access we require that no threads simultaneously update the same row in  $\hat{n}_{wt}$  matrix, yet the data for distinct words can be written in parallel. This is enforced by spin locks  $l_w$ , one per each word in the dictionary  $W$ . At the end of `ProcessDocument` we loop through all  $w \in d$ , acquire the corresponding lock  $l_w$ , append  $\hat{n}_{wt}$  to  $n_{wt}$  and

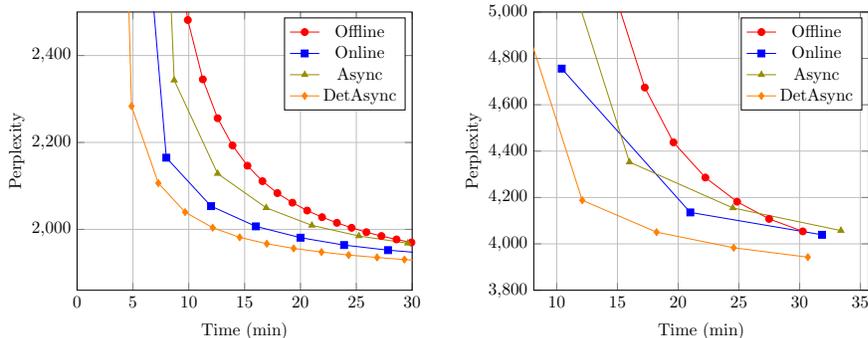


Fig. 8. Perplexity versus time for Pubmed (left) and Wikipedia (right),  $|T| = 100$  topics

Table 1. BigARTM peak memory usage, GB

	$ T $	Offline	Online	DetAsync	Async (v0.6)
Pubmed	1000	5.17	4.68	8.18	13.4
Pubmed	100	1.86	1.62	2.17	3.71
Wiki	1000	1.74	2.44	3.93	7.9
Wiki	100	0.54	0.53	0.83	1.28

release the lock. This approach is similar to [13], where the same pattern is used to update a shared stated in a distributed topic modeling architecture.

In our new architecture we also removed *DataLoader* thread, which previously was loading batches from disk. Now this happens directly from processor thread, which simplified the architecture without sacrificing performance.

In addition, we provided a cleaner API so now the users may use simple `FitOffline`, `FitOnline` methods to learn the model, and `Transform` to apply the model to the data. Previously the users had to interact with low-level building blocks, such as `ProcessBatches` routine.

## 8 Experiments

In this section we compare the effectiveness of `Offline` (Alg. 2), `Online` (Alg. 3), `Async` [17] and `DetAsync` (Alg. 4) algorithms. According to [17] `Async` algorithm runs approx. 10 times comparing to Gensim [10], and twice faster comparing to Vowpal Wabbit (VW) [11] in a single thread; an with multiple threads BigARTM wins even more.

In the experiments we use *Wikipedia* dataset ( $|D| = 3.7M$  articles,  $|W| = 100K$  words) and *Pubmed* dataset ( $|D| = 8.2M$  abstracts,  $|W| = 141K$  words). The experiments ran on Intel Xeon CPU E5-2650 v2 system with 2 processors, 16 physical cores in total (32 with hyper-threading).

Fig. 8 show the *perplexity* as a function of the time spend by the four algorithms listed above. The perplexity measure is defined as

$$\mathcal{P}(D, p) = \exp\left(-\frac{1}{n} \sum_{d \in D} \sum_{w \in d} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td}\right), \quad (5)$$

where  $n = \sum_d n_d$ . Lower perplexity means better result. Each point on the figures corresponds to a moment when the algorithm finishes a complete scan of the collection. Each algorithm was time-boxed to run for a 30 minutes.

Table 1 gives peak memory usage for  $|T| = 1000$  and  $|T| = 100$  topics model on Wikipedia and Pubmed datasets.

## 9 Conclusions

We presented a deterministic asynchronous (**DetAsync**) online algorithm for learning additively regularized topic models (**ARTM**). The algorithm supports all features of ARTM models, including multi-modality, ability to add custom regularizers and ability to combine regularizers. As a result, the algorithm allows the user to produce topic models with a rich set of desired properties. This differentiates ARTM from the existing models, such as LDA or PLSA, which give almost no control over the resulting topic model.

We provided an efficient implementation of the algorithm in **BigARTM** open-source library, and our solution runs an order of magnitude faster than the alternative open-source packages. Comparing to the previous implementation we eliminated certain performance bottlenecks, achieving optimal CPU utilization without requiring user to tweak batch size, number of inner loops per document, and other parameters. In addition, **DetAsync** algorithm guaranties a deterministic behavior, which makes it easier for us to unit-test our implementation and makes **BigARTM** ready for production use-cases.

In the future we will focus on memory efficiency to benefit from sparsity of word-topic ( $\Phi$ ) and topic-document ( $\Theta$ ) matrices, and extend our implementation to run on a cluster.

**Acknowledgements.** The work was supported by Russian Science Foundation (grant 15-18-00091). Also we would like to thank Prof. K. V. Vorontsov for constant attention to our research and detailed feedback to this paper.

## References

1. D. M. Blei. Probabilistic topic models. *Communications of the ACM*, 55(4):77–84, 2012.
2. Daud, A., Li, J., Zhou, L., Muhammad, F.: Knowledge discovery through directed probabilistic topic models: a survey. *Frontiers of Computer Science in China* 4(2), 280–301 (2010)

3. J. Yuan, F. Gao, Q. Ho, W. Dai, J. Wei, X. Zheng, E. P. Xing, T.Y. Liu, and W. Y. Ma. LightLDA: Big Topic Models on Modest Computer Clusters. In Proceedings of the 24th International Conference on World Wide Web, pp. 1351-1361, 2015.
4. D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
5. T. Hofmann. Probabilistic latent semantic indexing. In *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*, pages 50–57, 1999.
6. M. D. Hoffman, D. M. Blei, and F. R. Bach. Online learning for latent dirichlet allocation. In *NIPS*, pages 856–864. Curran Associates, Inc., 2010.
7. D. Newman, A. Asuncion, P. Smyth, and M. Welling. Distributed algorithms for topic models. *J. Mach. Learn. Res.*, 10:1801–1828, Dec. 2009.
8. Y. Wang, H. Bai, M. Stanton, W.-Y. Chen, and E. Y. Chang. PLDA: Parallel latent Dirichlet allocation for large-scale applications. In *Proceedings of the 5th International Conference on Algorithmic Aspects in Information and Management*, pages 301–314, 2009.
9. Z. Liu, Y. Zhang, E. Y. Chang, and M. Sun. PLDA+: parallel latent Dirichlet allocation with data placement and pipeline processing. *ACM Trans. Intell. Syst. Technol.*, 2(3):26:1–26:18, May 2011.
10. R. Řehůřek and P. Sojka. Software framework for topic modelling with large corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta, May 2010.
11. J. Langford, L. Li, and A. Strehl. Vowpal wabbit open source project. *Technical report*, Yahoo!, 2007.
12. A. K. McCallum, A Machine Learning for Language Toolkit. <http://mallet.cs.umass.edu>, 2002.
13. A. Smola and S. Narayanamurthy. An architecture for parallel topic models. *Proc. VLDB Endow.*, 3(1-2):703–710, Sept. 2010.
14. K. V. Vorontsov. Additive regularization for topic models of text collections. *Doklady Mathematics*, 89(3):301–304, 2014.
15. K. V. Vorontsov and A. A. Potapenko. Tutorial on probabilistic topic modeling: Additive regularization for stochastic matrix factorization. In *AIST’2014, Analysis of Images, Social networks and Texts*, volume 436, pages 29–46. Springer International Publishing Switzerland, Communications in Computer and Information Science (CCIS), 2014.
16. K. V. Vorontsov and A. A. Potapenko. Additive regularization of topic models. *Machine Learning, Special Issue on Data Analysis and Intelligent Optimization*, volume 101(1), pages 303–323, 2015.
17. K. Vorontsov, O. Frei, M. Apishev, P. Romov and M. Dudarenko. BigARTM: Open Source Library for Regularized Multimodal Topic Modeling of Large Collections. In *AIST’2015, Analysis of Images, Social networks and Texts*, volume 542, pages 370–381. Springer International Publishing Switzerland, Communications in Computer and Information Science (CCIS), 2015.
18. K. Vorontsov, O. Frei, M. Apishev, P. Romov, M. Suvorova, A. Yanina. Non-bayesian additive regularization for multimodal topic modeling of large collections. In *Proceedings of the 2015 Workshop on Topic Models: Post-Processing and Applications*, pp. 29–37. ACM, New York, USA, 2015.