

Two Independent and Highly Efficient Open Source TKF91 Implementations

Nikolai Baudis, Pierre Barbera, Sebastian Graf, Sarah Lutteropp, Daniel Opitz, Tomas Flouri, and Alexandros Stamatakis

Karlsruhe Institute of Technology, Institute of Theoretical Informatics,
Kaiserstrasse. 12, 76131 Karlsruhe, Germany
{daniel.opitz,sebastian.graf,nikolai.baudis}@student.kit.edu
<http://www.informatik.kit.edu/>

Abstract. In the context of a master level programming practical at the computer science department of the Karlsruhe Institute of Technology, we developed and make available two independent and highly optimized open-source implementations for the pair-wise statistical alignment model, also known as TKF91, that was developed by Thorne, Kishino, and Felsenstein in 1991. This paper has two parts. In the educational part, we cover teaching issues regarding the setup of the course and the practical and summarize student and teacher experiences. In the scientific part, the two student teams (Team I: Nikolai, Sebastian, Daniel; Team II: Sarah, Pierre) present their solutions for implementing efficient and numerically stable implementations of the TKF91 algorithm. The two teams worked independently on implementing the same algorithm. Hence, since the implementations yield identical results—with slight numerical deviations—we are confident that the implementations are correct. We describe the optimizations applied and make them available as open-source codes in the hope that our findings and software will be useful to the community.

1 Introduction

In [1], Thorne, Kishino, and Felsenstein presented a method for aligning DNA sequences using a maximum likelihood approach. To this end, they developed an explicit statistical model of evolution that includes insertions, deletions, and substitutions of nucleotides as basic operations for comparing two DNA sequences. A given evolutionary model assigns a weight to each of these three basic operations. These are then used to design a dynamic programming algorithm that computes the maximum likelihood pair-wise alignment between two DNA sequences.

As substitution model, we assume the standard stochastic nucleotide substitution model **TODO by whom actually?. TODO the other two?** The algorithm computes the optimal maximum likelihood sequence alignment using three matrices, one for each **TODO**. The algorithm is given by the following

dynamic programming recurrence:

$$\begin{aligned}
M^0(i, j) &= \frac{\lambda}{\mu} \pi_{a_i} \overline{p_0}(t) \max\{M^0(i-1, j), M^1(i-1, j), M^2(i-1, j)\} \\
M^1(i, j) &= \frac{\lambda}{\mu} \pi_{a_i} \max\{P_{a_i \rightarrow b_j}(t) p_1(t), \pi_{b_j} \overline{p_1}(t)\} \\
&\quad \max\{M^0(i-1, j-1), M^1(i-1, j-1), M^2(i-1, j-1)\} \\
M^2(i, j) &= \pi_{b_j} \lambda \beta(t) \max\{M^1(i, j-1), M^2(i, j-1)\}
\end{aligned}$$

Where λ and μ denote the birth and death rate, $\pi_x, x \in \{A, C, G, T\}$ denotes the equilibrium frequency of the nucleotides, $P_{x \rightarrow y}$ the transition probability from x to y , and $\overline{p_n}(t)$ the probability that after time t a mortal link has exactly n descendants. The values for t , λ , μ , and π are given as input. The values for $\overline{p_n}$ can be pre-computed in constant time **constant?**.

TODO related work

TODO acronyms: DP, DPM

The remainder of this paper is organized as follows: In Section 2, we describe the teaching setup and goals. The teams then present their implementations in Sections 3 and ??, respectively. The corresponding experimental results by both teams are presented in Section 5. In the following Section 6, we summarize our teaching experiences. We conclude in Section 7.

2 Teaching Perspective, Goals and Course Outline

2.1 Teaching Setup & Goals

Courses at the Master level in our computer science department are organized in so-called modules over two semesters. In the first semester of the Bioinformatics module, we teach a lecture called “Introduction to Bioinformatics for Computer Scientists”, since KIT does not offer a stand-alone Bioinformatics degree. This lecture covers basic topics such as an introduction to molecular biology, pairwise sequence alignment, BLAST, de novo and by-reference sequence assembly, multiple sequence alignment, phylogenetic inference, MCMC methods, and population genetics.

In the second semester of the module, students can choose if they want to do a seminar presentation or the programming practical whose results we describe here. The goal of the practical is to carry out a self-contained project and write, as well as release software, that is useful to the evolutionary biology community. Another focus is on learning to use tools that enhance software quality. Note that, at a CS department designing “classic” bioinformatics analysis pipelines using scripting languages is typically not considered as “real programming” by the students. Hence, we needed to define a project that requires coding in **C/C++** or **Java**. One should also strive to avoid having the students extend existing software, since this is generally frustrating and hinders creativity.

We thus decided to ask the students to implement efficient, sequential versions of the TKF91 algorithm that can also be used as library routines. Since the

dynamic programming algorithm as such, is relatively straight-forward to implement for a computer science student, the main focus was on code optimization. The students were thus asked to implement a highly efficient version of the code in C or C++ using all capabilities of a modern CPU (e.g., SSE3 and AVX intrinsics). The main motivation for this was to give the students enough time to experiment with low-level optimization strategies on modern CPUs. Moreover, this project allowed to apply a broad range of skills acquired in the Bioinformatics and other master-level modules at our department. Furthermore, the TKF91 model required understanding and applying the discrete pair-wise sequence alignment methods and likelihood-based models for evolutionary biology presented in the lectures. To foster competition among the teams, an award (dinner payed by the PI) was announced for the team that would implement the fastest code. To allow for a fair comparison of the codes and CPU-specific optimization we provided the students access to a reference machine with 4 physical cores (Intel i7-2600 running at 3.40GHz) and 16GB RAM.

In terms of project documentation, students are usually required to write a report. However, in the present case, we jointly took the decision to write a paper about the practical and upload it to [biorxiv](https://www.biorxiv.org/). This has the positive effect that students also learn how to write scientific papers.

2.2 Code Quality Assessment

In order to continuously monitor and improve the quality of our source code, we used several tools and methods throughout our project. We deployed both, static, as well as dynamic code analyses.

Static Analysis Static analyses help to identify static programming errors such as incorrect programming language syntax or static typing errors. We compiled our code with the `gcc` compiler using all available and reasonable warning flags¹. This allowed us to identify potential programming errors at an early stage. Due to its more pedantic nature (i.e., ability to detect more errors), we also used the `clang` compiler with respective flags² periodically alongside of `gcc` to further reduce the amount of potential errors.

Dynamic Analysis These analyses cover run-time issues such as memory leaks. We used `valgrind` and its sub-module `memcheck` for detecting memory-related errors in our programs.

¹ -Wall -Wextra -Wredundant-decls -Wswitch-default -Wimport -Wno-int-to-pointer-cast -Wbad-function-cast -Wmissing-declarations -Wmissing-prototypes -Wnested-externs -Wstrict-prototypes -Wformat-nonliteral -Wundef

² -Weverything -pedantic

3 Implementation of Team I

Before describing our vectorized implementation, we cover the necessary prerequisites. In Section 3.1 we describe the DP algorithm in more detail as well as the corresponding wave-front parallelism we exploited in our implementation. In Section 3.2, we describe our memory layout concept, which is required to improve the efficiency of the vectorization. Thereafter, in Section 3.3, we describe how we cache intermediate computations, and then explain our vectorization approach in Section 3.4. In addition to that, Section 3.5 outlines our considerations regarding the alignment of the data, marking an important part of the vectorized implementation.

3.1 Dynamic Programming and Wave-front Parallelism

Dynamic Programming (DP) is a widely-used technique to efficiently solve a class of, at first sight, apparently hard computational problems when they exhibit *overlapping subproblems* and *optimal substructure*. Pair-wise sequence alignment falls into this class. Another advantage of DP is that it is comparatively straightforward to parallelize computations along the anti-diagonals of the DP matrix. This DP matrix parallelization approach is also known as wave-front parallelism. The underlying idea is that matrix entries along the same DP anti-diagonal d can be computed independently from each other (and hence in parallel) if the preceding anti-diagonal $d-1$ has been computed. Therefore, we can deploy vector intrinsics to accelerate calculations along anti-diagonals.

3.2 Memory Layout

We first introduce our memory layout for the three DP matrices (M^0, M^1, M^2) since this is performance-critical. Our vectorization approach needs to attain high data locality to take advantage of the CPU cache. To this end, we permuted the indexing scheme of the DP matrix such that matrix neighboring entries on a anti-diagonal are also contiguous in memory. Figure 1 shows the memory mapping of the matrix. It is layed out neither in a column-major nor row-major fashion in memory, but stored linearly by anti-diagonals.

We index the matrices using a modified version of the *Cantor pairing function*, where a tuple is assigned to an offset. The original Cantor pairing function π assigns an integer to a pair of integers (e.g., $\pi : (\mathbb{N}, \mathbb{N}) \rightarrow \mathbb{N}$). Since the size of the DPM is defined by the length of the input sequences and because it is not infinite, we modified the pairing function to accept a tuple of two integers from the interval $[0..length(\text{Sequence } 1|2)]$ as input. As a consequence, the DPM is divided into three parts: the *opening*, *middle*, and *closing part*. In the opening part, the modified pairing function is identical to the original function; in the middle and closing parts, the modified pairing function is calculated by subtracting an offset from the original pairing function. We calculate the offset via the original Cantor pairing function.

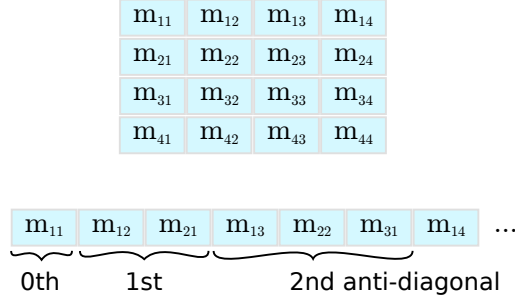


Fig. 1. Memory layout for the DP matrix

0	1	3	6-0	10-1	0	1	3
2	4	7-0	11-1	16-3-1	2	4	
5	8-0	12-1	17-3-1	23-6-3	5		
0	1	3	6				
2	4	7					
5	8						

Fig. 2. Indexing scheme for the anti-diagonals

In Figure 2, we show the index calculations for the anti-diagonals of a 5×3 -matrix. The blue box denotes the boundaries of the DP matrix, the blue cells represent the opening part, the green cells inside the matrix the middle, and the red cells inside the matrix the closing part. Green and red cells outside of the matrix boundaries represent the offset calculation for the middle and closing part, respectively. For the DP calculations **for what?**, only the indices that correspond to entries in the DP matrix **of those all? pairing functions** are needed.

Given the anti-diagonal indexing scheme for a single matrix, we need to devise the memory layout for the three DP matrices (M_0, M_1, M_2) next. This is because calculating a DP value requires accessing values from all three matrices, which are located at different positions in memory. **I am not sure I understand the following question, please re-formulate – Simplified the sentences, the Therefore-sentence wasn't really required i think.** To improve data locality while, at the same time, using efficient **vector?** vector load and store operations, the matrix data has to be arranged accordingly. Since we are using double-precision floating-point numbers for all calculations, an SSE3 vector (128 bit) holds 2 values, while an AVX vector (256 bit) holds 4 values. To this end, we evaluated the following three alternative matrix layouts (see Figure 3):

Struct of Arrays (SoA) (`struct { double m0[n], m1[n], m2[n]; } data;`) This layout allows for simple **vector?** vector load and store patterns. However, the

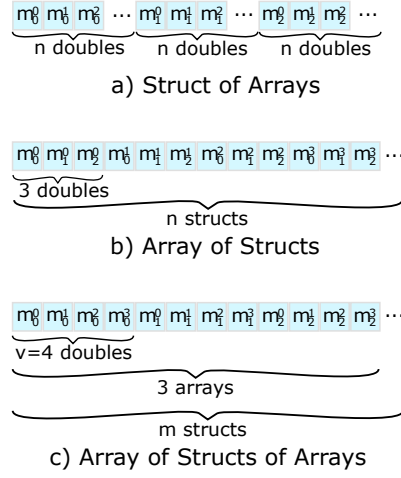


Fig. 3. Alternative data layouts for storing anti-diagonals

three matrices are located in separate memory regions, which decreases cache efficiency and data locality.

Array of Structs (AoS) (`struct { double m0, m1, m2; } data[n];`) This layout exhibits improved data locality by grouping values closely together that are – in most cases – accessed simultaneously. On the other hand, loading and storing vectors requires disentangling and interleaving them, which requires several, potentially costly, vector shuffle operations.

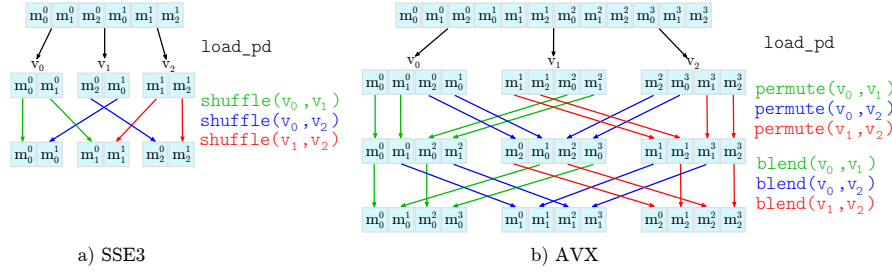


Fig. 4. Disentangling operations for a) SSE3 (128 bit) and b) AVX (256 bit) vectors.

Figure 4 illustrates these shuffle operations. For m_x^y , the subscript x indicates the index of the matrix and the superscript y the SSE3/AVX vector element index y **state that we use double precision! Added a sentence about double-precision and vector size above.** The first row of the Figure depicts the *AoS* data layout. Here, we perform a **if you say aligned you need to be**

more specific– Removed the "aligned", I think it's not important while discussing the shuffle operations vector load operation to load the data into temporary vector registers. Thereafter, we shuffle the temporary vectors such that (i) each vector only contains elements from one of the tree DP matrices and (ii) the individual elements are in the correct order with respect to the **anti-?** anti-diagonal. For both SSE3 as well as AVX instructions, we need to perform three load operations and three **shuffle/permute** operations. For AVX, three additional **blend** operations are required to put the data elements into the correct order.

Array of Structs of Arrays (AoSoA) (`struct { double m0[v], m1[v], m2[v]; } data[m];`) This layout groups vectors in advance, i.e., adjacent elements of each matrix are already in an SSE3/AVX vector. While solving both of the previous problems **state those problems again!** (low data locality and costly vector operations), grouping vectors **what do you mean by grouping vectors?– Added an explanation at the beginning** in advance complicates DP cell access above the current anti-diagonal. These values do **would or do not lie?** not lie within vector boundaries and would **would or do require?** then again **either?** require either **shuffle/permute** operations or partial loads **explain what you mean by partial loads!** (i.e., accessing single elements of several vectors).

All three layouts have different performance characteristics with respect to cache efficiency **Note to team I: we usually say cache efficiency and data locality, but not cache locality!** and complexity of the load, store, and shuffle operations. We implemented prototypes for all three approaches to analyze their performance and found that the AoS approach performs best. Note that, the complexity of the **shuffle** operations outweighed the improved cache efficiency in the AoS and AoSoA layouts (also see Table 5.2 in Section ??).

3.3 Memoization

The [1] model has several parameters. Thus, it initially seemed inevitable to carry out the non-trivial computations for filling DP cells for each cell from scratch. However, after analyzing which expressions are constant and can thus be pre-computed and reused (i.e., memoized **Note: this could also be called caching or using a lookup table**), we reduced the operations required for calculating a DP cell value to just three: indirection, summation, maximum.

Obvious savings can be achieved for expressions that only depend on the time parameter t , the birth rate λ , and the death rate μ . These values are constant and given as input parameters. We also observed that the cell updates in the DP algorithm only require a constant number of common sub-calculations/components: When the nucleotide values at the current indices in the two sequences are available, the score (cell value) can be computed using only the current nucleotide pair and the neighboring values **of which matrices** of the three DP matrices.

The memoization scheme (lookup table) for all possible transition penalties is shown in Figure 5. The penalties C^i take one or two nucleotides as parameters. Thus they can be stored in a memoization table of 4 and 16 entries, respectively. With a total of only 24 memoized values, we substantially simplified and accelerated (by **TODO how many orders of magnitude? Actually, I can't tell, because we hadn't had our benchmark suite in place back then. I could cherry pick some commits and throw together a benchmark run if that's really wanted, but I'd rather get rid of that claim.**) the cell updates. In addition, this simplification now allows to apply the logarithm for preventing numerical underflow. This was not possible before, since taking the logarithm of the original equation would have been too expensive. It also kept the translation to SIMD-enabled code manageable **what do you mean by manageable? The complexity involved with having to translate the recurrence steps verbatim to SIMD instructions and optimizing them with confidence. I imagined that would be much more involved and harder to parallelize efficiently.** Furthermore, the simplification reduces the number of design decisions that are required regarding memory layout and vectorization.

$$\begin{aligned}
C^0(a) &= \frac{\lambda}{\mu} \pi_a \overline{p_0}(t) \\
C^1(a, b) &= \frac{\lambda}{\mu} \pi_a \max\{P_{a \rightarrow b}(t) p_1(t), \pi_b \overline{p_1}(t)\} \\
C^2(b) &= \pi_b \lambda \beta(t) \\
M^0(i, j) &= C^0(a_i) \max\{M^0(i-1, j), M^1(i-1, j), M^2(i-1, j)\} \\
M^1(i, j) &= C^1(a_i, b_j) \max\{M^0(i-1, j-1), M^1(i-1, j-1), M^2(i-1, j-1)\} \\
M^2(i, j) &= C^2(b_j) \max\{M^1(i, j-1), M^2(i, j-1)\}
\end{aligned}$$

Fig. 5. A re-formulation of the DP step using the memoized sub-problems C^i **you should have explained what the C is much earlier! I think we referenced it much earlier, so I moved this nearer to its 'declaration site'.**

At a later point of the project, it became evident that the indexed loads from the memoized penalty matrices C^0 , C^1 , C^2 as defined in Figure 5 caused a performance degradation. To solve this problem, we tried to pre-compute a larger lookup table of vector-sized elements, that we indexed by a specific nucleotide permutation. Consider the following example for AVX vector intrinsics. We need to calculate a bijective mapping for a set of four nucleotides to an integer representing one of the $4^4 = 256$ possible permutations (e.g. a perfect hash) for C^0 and C^2 , and $4^{2 \cdot 4} = 65536$ for C^1 respectively **where do those C variables come from, have they been properly introduced? I don't think so! Those are the cached values that were on the next page. Moved them**

up and added an introductory phrase at the begin of the paragraph.. The mapping is implemented as a base conversion from the set of 4 nucleotides to an unsigned integer via appropriate bit-wise operations.

While this approach has exponential space requirements as a function of the vector width, using a lookup table of $65536 \cdot 32$ bytes = 2 MB for C^1 for the AVX version of our code was still feasible. However, the evaluation of the resulting performance revealed that too much time is spent to fill the table, at least for our benchmark datasets. Since the values in the table only depend on the likelihood model parameters, that is, it is independent of the input sequences at hand, this approach can be re-considered in cases were this lookup table shall be re-used for multiple DP invocations. This might be useful when calculating pair-wise alignments of n sequences under the same model parameters **not sure if we should write this here, at leats as long as the parameter t is also required to calculate the values in that table, need to double-check**That's a good point. I was mentally pulling t into the model, which is incorrect, unless t stays constant. I don't understand the following sentence, please re-formulate it Tried to clarify with the part following e.g. The scaling to wider vector registers could be retained by loading a register in packs of four values, e.g. for an 8 element vector, loading the lower and upper 4 values separately from the lookup table for 4 element vectors allows the space complexity to stay constant with respect to vector width.

3.4 Vectorization

We implemented a vectorized version for computing M_0 , M_1 , and M_2 using `add`, `max`, `load`, and `store` operations for both SSE3 and AVX instructions. Using these operations, we can process n matrix elements in parallel. Depending on the location of the vector on the anti-diagonal that is being processed, there will be $0 < n \leq V$ valid elements to operate on per vector (with vector size $V = 4$ for AVX and $V = 2$ for SSE3 **TODO: we need to state in the beginning that we used double precision**Added a sentence about double-precision and vector size above).

Usually, data that has to be loaded into vectors is aligned, i.e., the starting address of the data is within certain byte boundaries. This alignment allows to use the aligned versions of the vector `load` and `store` operations, which are faster than the unaligned versions. Although we ensured aligned **need to be more specific what you mean by aligned, it's not clear to a general readership**Added a sentence about alignment access to a majority of the data (as will be explained in Section 3.5), we used the unaligned vector `load` and `store` operations `loadu_pd` and `storeu_pd` for the SSE3 and AVX versions. **I don't get the following sentence, whyt was unaligned load necessary if you were able to test on aligned data???**– We aligned the data so that the majority of the vectors start at an aligned address. If, for example, $M_2(i-1,j)$ (the same position shifted by -1 vertically) needs to be loaded into a vector, we can only do that using an unaligned load. This does not represent a disadvantage for the performance, since we could not

identify significant performance impacts of using unaligned intrinsics on aligned data.

To prevent our `store` operations from overwriting data at the boundaries of the matrices, we ensured that for vectors with size $n < V$ (i.e., not completely filled/padded vectors) only n elements are written back into the matrices. For the SSE3 version, this can be achieved by only writing the lower part of the vector if its size is 1. For AVX, however, we covered all cases where $n < V$ holds using the `maskstore_pd` intrinsic and an appropriate mask for all possible value of n (i.e., 3, 2, and 1).

3.5 Data Alignment

When vectorizing along the anti-diagonals, accessing the memory in the dynamic programming matrices when they are stored in a row-, or column-major scheme prevents deploying efficient vector operations for loading the values from the matrices into vector registers. For such a matrix layout, the performance of the inner DP loop becomes heavily memory-bound. Therefore, we used the aforementioned diagonalized matrix layout together with an struct of arrays approach (see section 3.2) **TODO: add cross-ref done**, such that `load` and `store` operations for copying data directly from the anti-diagonal into vector registers can be utilized. Those `load` and `store` operations come in two flavors: There exist `load` and `store` operations for aligned and unaligned data. Using the unaligned versions can decrease memory access performance. However, as mentioned in Section 3.4, loading aligned data using `loadu_pd` does not incur a performance penalty. On the other hand, using aligned loads yields storing the DP matrices more complex. When loading data from the matrices into SSE3/AVX vector registers, the data has to be aligned to 16/32 byte boundaries.

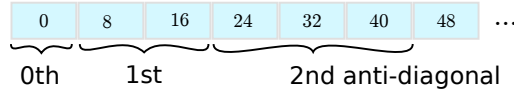


Fig. 6. Shift of the byte offset of the anti-diagonals

Figure 6 shows how the byte offset for the anti-diagonals is shifted from one diagonal to the next **to achieve a proper 16/32 byte alignment I assume? Not quite, it's how alignment to 16/32 byte boundaries is not guaranteed when just tightly packing anti-diagonal after anti-diagonal linearly into the buffer. We modified the indexing function so that it inserts the appropriate padding cells, such that the beginning (well, actually each 2nd cell, the first was initialized sequentially before the tight loop and thus can be unaligned) of an anti-diagonal is always aligned. Check out the comment here: <https://goo.gl/r8H6qF>. It's pretty technical, but I hope this clears things up. Note how the entire first row is unaligned and how the second row is always aligned**

to multiples of 4. To easily determine the aligned addresses for loading vectors we modified the indexing function for diagonalized access (see also section 3.2) **TODO cross ref! done** for the anti-diagonals. The main idea is to store each anti-diagonal starting at a 16/32 byte-aligned address. **is the following edit I mad actually correct? It's not that all load and stores are aligned (that's not possible), but that we choose a layout where the most possible loads/stores are aligned.** Thereby, all load and store operations are guaranteed to happen on a predictable alignment, which we choose so that most of the operations happened on a proper alignment boundary.

The following three sentences need to be re-written, they are rather confusing I rewrote them. Is it clearer now? We strived for a scheme where the row with index 1 (zero-based) is aligned, because the row with index 0 is initialized prior to the DP loop. This resulted in each anti-diagonal to to be aligned on a rather strange alignment ($a - 8$ where a is the alignment boundary in bytes, so 24 bytes for the 32 byte example).

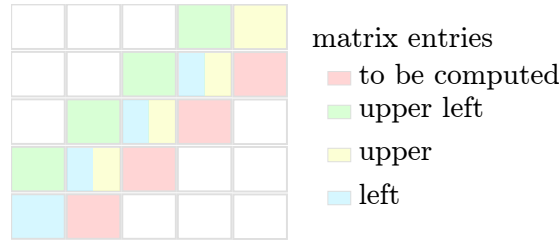


Fig. 7. Outline of the data required for an operation in the DP loop

In Figure 7, we show the required values for computing the inner DP loop for AVX (vector length: 4 double precision values). The elements marked in red are the ones we want to compute in the current step. For calculating a single element, we require three elements from the two previous anti-diagonals. The element directly above (yellow) and left (blue) to the current element and the diagonal element above and to the left (green) need to be loaded from the DP matrices. The Figure illustrates that (neglecting the asymptotically irrelevant boundary cases) one of the following four conditions holds:

1. red and blue start on an aligned address, yellow and green do not
2. red and yellow start on an aligned address, blue and green do not
3. yellow and green start on an aligned address, red and blue do not
4. red and green start on an aligned address, yellow and blue do not

For the computation, we need to perform three **load** operations in total for the blue, yellow, and green elements and one **store** operation for the red elements. Since the *unaligned store* operation has greater relative performance impact than the unaligned load operations **why?Good question :) purely speculative, I guess. Although stores seem to take longer than loads.**

, we ensured that either the first or the second of the above conditions is met. The first condition is implemented for the opening and middle part, the second condition holds true for the closing part of the DP matrices.

We achieved an appropriate alignment for each anti-diagonal by padding anti-diagonals as needed.

4 Implementation of Team II

In the following, we first describe how we transformed the algorithm into log-space to prevent numerical underflow. This also allowed us to simplify the formulas. Thereafter, we report how we improved data locality by storing the matrix entries in a dedicated data structure. While we experimented with different vectorization techniques, it turned out that the fastest code did not rely on vectorization. We thus conclude that, we managed to simplify the sequential code to such an extent, that, the potential advantages of using SIMD instructions are amortized by the additional overhead (**TODO add cross-ref to the section where you discuss this**) we need for vectorization.

4.1 Mathematical Optimization

Numerical Underflow Prevention As the TKF91 algorithm performs successive multiplications on floating point numbers, preventing numerical underflow *is* a major issue. Underflow can occur, even for short input sequences with less than 100 nucleotides each. To address this issue we transformed all computations into log-space, to add logarithms of probabilities instead of multiplying probabilities. While it is still possible to experience under- or overflow **why overflow, please explain! Pierre: I think we meant when transforming a problem into logspace in general, but reading on the text is very specific, so I would remove the “overflow”** even after this transformation, this is highly unlikely to occur for practical (empirical) input data. We did not observe any numerical issues for the range of input parameters values and sequence lengths (up to and including 10000 nucleotides) we tested.

Simplifying the Formulas We were able to omit redundant computations by re-using already computed values from previous matrix entries. For example, we observed that

$$M^0(i+1, 0) = M^0(i, 0) + \log(\gamma_{i+1}) + \log(\zeta_{i+1}) + \log(\beta(t)) + \log(\pi_{a_{i+1}}) + \log(\bar{p}_0(t)).$$

After replacing $\beta(t)$, $\bar{p}_0(t)$, γ_i and ζ_i by their respective formulas, we observed that some terms appear multiple times. Operating in log-space allowed us to further simplify the formulas. Especially the logarithmic rules $\log(a*b) = \log(a) + \log(b)$ and $\log(a/b) = \log(a) - \log(b)$ allowed us to replace multiplications and divisions with additions and subtractions. The resulting formulas can be found in Appendix ?? **please include them here!**

Pre-computing the Logarithms In general, replacing a multiplication by two logarithms and an addition is more costly. However, we managed to circumvent the additional computational cost of log space. This is because our formulas consist of sums of only 25 different logarithms **explain why this is the case** which we can simply pre-compute.

A major disadvantage of logarithmic transformations is that any error in the logarithm computation accumulates over a series of operations. To further investigate this, we used the high precision mathematical library of the `boost`-framework [?]. It offers datatypes that dynamically adapt their floating point precision to avoid numerical errors as well as under-/overflow. Our initial idea was to use this library only for pre-computing the logarithms, and thereby reduce the induced runtime overhead. However, transforming these arbitrary precision datatypes to standard double precision floating point values proved to be problematic, since it introduced additional loss of precision **comapred to what?**. Therefore, we abandoned this path.

4.2 Matrix Storage Schemes

Matrix as Array In the first, naïve implementation, we stored each matrix (M^0, M^1, M^2) in a separate array, using row-major order. The matrix entry at position (i, j) is stored at index position $i * (m + 1) + j$ in the array (see Figure 10).

To illustrate the shortcomings of this approach, consider the following example: assume that six matrix rows fit into one cache line. Then, performing one iteration of the inner DP loop requires loading three cache lines, one per matrix. This is depicted in part a) of Figure 8. If we further assume a cache capacity of only two cache lines, every iteration would then force one of the lines to be swapped out, causing memory overhead.

Array-of-Structs Data Structure The dynamic programming step (see Algorithm 1) accesses the matrices M^0, M^1 and M^2 at the same index position to determine the maximum entry. Thus, we improved data locality by storing the entries from the three matrices at the same index position contiguously in memory. For this, we implemented a data structure called `MatrixEntry` that consists of three double values: m_0, m_1 and m_2 . Then, we used the Matrix as Array approach as before, but containing `MatrixEntry` structs as elements instead of simple `double` values (see Figure 9).

Let us now consider the example in Figure 8 again. With six matrix lines filling one cache line, a cache line now contains data from all three matrices. Consequently, all operations for a single computational iteration **what do you mena by computational iteration?** will need to access two cache lines at most. If we assume a cache capacity of two cache lines again, cache misses will now only occur when traversing row boundaries.

By using the `perf stat` tool, we found that using the `MatrixEntry` data structure indeed reduced **by how much, please quantify!** the number of page faults **is this really called page faults? shouldn't it be cache misses?**.

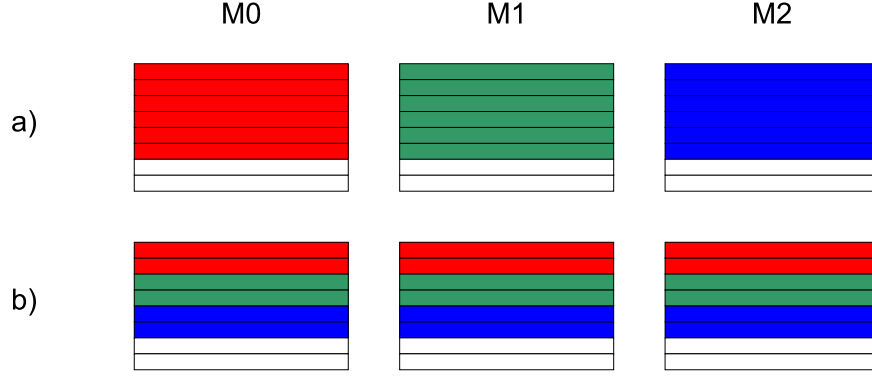


Fig. 8. Relation of cache lines and matrices, using the method of allocating each matrix separately **(a)** and the Array-of-Structs method **(b)**. Each color represents a distinct cache line.

Algorithm 1: The dynamic programming step, row-major version

```

1 ...// some initializations, see Appendix ??
2 for  $i = 1, \dots, n$  do
3   for  $j = 1, \dots, m$  do
4     ...// some initializations, see Appendix ??
5      $coord \leftarrow CO(i, j)$ ;
6      $up \leftarrow CO(i, j - 1)$ ;
7      $diag \leftarrow CO(i - 1, j - 1)$ ;
8      $left \leftarrow CO(i - 1, j)$ ;
9      $m[coord].m_0 \leftarrow m[coord].m_0 + \max\{m[left].m_0, m[left].m_1, m[left].m_2\}$ ;
10     $m[coord].m_1 \leftarrow$ 
11     $m[coord].m_1 + \max\{m[diag].m_0, m[diag].m_1, m[diag].m_2\}$ ;
12     $m[coord].m_2 \leftarrow m[coord].m_2 + \max\{m[up].m_1, m[up].m_2\}$ ;
13  end
14 end

```

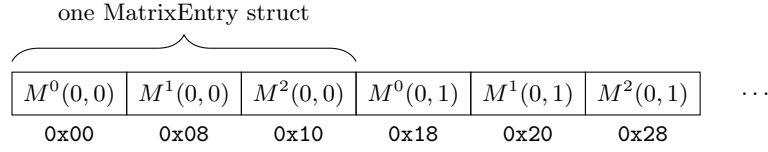


Fig. 9. Array-of-Structs data structure used to store the three matrices in memory. Each `MatrixEntry` element stores the entries of a single coordinate for all three matrices. The structs are stored contiguously, as the figure depicts, in row-major order.

Alternative Storage Schemes TODO: need to think how to match/connect this with the indexing that team I developed

Since each iteration in the dynamic programming step needs to access the top, left, and upper-diagonal elements of the matrices, a wave-front-like storage of the matrix (see Figure 11) entries yields the lowest amount of expected cache misses. Finding an easy-to-compute **do you mean easy or inexpensive?** closed formula that maps the index position (i, j) representing the row and column of a matrix to wave-front-coordinates turned out to be challenging. The index of the diagonal is $i + j$, the sum of row index and column index. Determining the number of elements in the diagonal and especially the correct position within the diagonal proved more difficult though. Since we were not able to come up with an easy **still not sure what you mean by easy** formula for the mapping, we tried storing pre-computed index mappings in an additional matrix. However, this slowed down the program as computing the mapped indices required more arithmetic operations than needed for the actual TKF algorithm. Figure 12 illustrates the main problem which has to be solved for effective wave-front indexing effective: Given an wave-front index k , how do we obtain the indices for the upper, upper-left-diagonal and left element of the matrix? While the required offsets are constant within a **anti?**diagonal, they change for different **anti?**diagonals in the matrix.

0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19

Fig. 10. Row-major indexing

0	1	3	6	10
2	4	7	11	14
5	8	12	15	17
9	13	16	18	19

Fig. 11. Wave-front indexing

4.3 Vectorization Attempts

During code development we conducted a partial vectorization for two of our implementations: the basic Log-space implementation (Section 4.1) and the im-

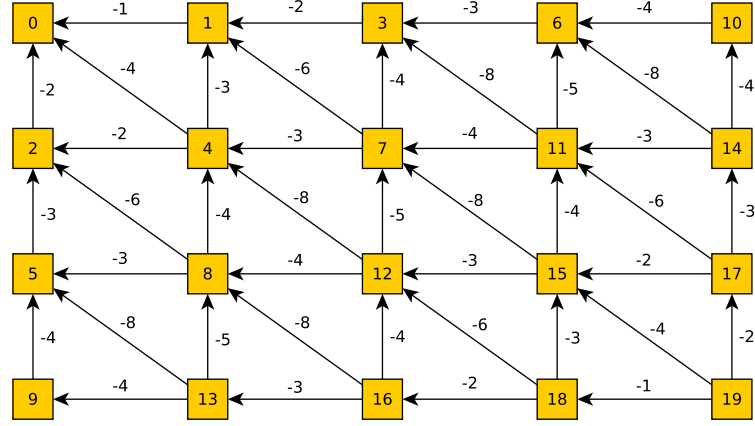


Fig. 12. Offsets in wave-front indexing

proved Log-space implementation, using a more cache-efficient data structure to store the matrices (Section 4.2).

In both cases we mainly focused on vectorizing the code that initializes the three matrices. This is because, our analyses with `Valgrind` `callgrind` revealed that this part of the code required 75 – 80% of overall runtime. Additionally, as the initializations are completely independent of each other, this part should be straight-forward to vectorize.

Several code transformations were necessary to vectorize the code. First, we needed to guarantee that vector `load` and `store` operations are invoked on correctly aligned memory addresses. Secondly, we needed to devise a strategy for iterating over matrices and handling matrix sizes that are not multiples of the vector width. Finally, we needed to store calculation terms, such as the $\log(\pi)$'s and $P_{a_i \rightarrow b_j}$'s, **what do you mean by calculation parameters?** **Pierre: added some clarification** into the appropriate vector types and had to replace all basic arithmetic operations by their vectorized counterparts.

Log-space implementation For the vectorization of the basic **did you properly introduce what the basic implementation is?** **Pierre: First paragraph in “Vectorization attempts” we name it as such, provide a back reference to 4.1. Maybe we could introduce it directly under the 4.1 section header? then back ref here again** Log-space implementation we implemented appropriate memory alignment by using the `__attribute__((aligned(size)))` attribute for stack allocations, and the `posix_memalign` function for heap allocations.

As a consequence, we also changed the iterations through the matrices to start at entry 0 of each row, as starting iteration at position 1 would result in incorrectly aligned memory accesses. Furthermore, as we stored the rows of the

matrices contiguously in memory, we had to pad the row sizes to ensure the correct alignment of the first entry in each row.

The strategy we chose to deal with row dimensions that are not multiples of the vector width is straight forward. Per row, we compute as much as possible via vector intrinsics and the remainder sequentially.

Version using Array of Structs data structure While the strategy for dealing with row dimensions that do not fit vector widths remained the same in the improved version of our program, we had to apply several modifications to use vector intrinsics in conjunction with the more cache-efficient data structure.

As Figure 9 shows, using the Array-of-Structs scheme results in a non-contiguous storage of row entries. As a consequence, simply loading and storing vector registers cannot be used here. To this end, we implemented load and store operations that fit our data structure. These operations rely on vector functions that load and store single double precision floating point values. An advantage of this is that correct data alignment is not required any more. However, we believe **why do you believe this, don't you have experimental evidence? Pierre: we didn't test them in isolation, but we assume that the load/store intrinsics that require contig. memory (`_mm_load_pd`) work faster than my function using `_mm_loadl_pd` and `_mm_loadh_pd`** that this approach also induces a significant performance overhead compared to load/store intrinsics that work on contiguous memory. **what are classically used operations? Pierre: changed it.**

Thoughts on wave-front vectorization Maybe somehow make the wave-front parts more consistent to avoid redundnacies? As the dynamic programming step of the TKF91 algorithm has data dependencies to the cells left of, above and diagonally left above of the cell to be computed, parallelization/vectorization requires a *wave-front* approach.

In a wave-front parallelization scheme, matrix entries are not computed by row or by column, but rather by computing the entries of each *anti-diagonal* in parallel. Crucially, the individual anti-diagonals have to be computed sequentially and in-order when dependencies exist between subsequent anti-diagonals. In this case, the current anti-diagonal is the front of the wave traversing the DP matrix.

The following sentence is too long and confusing, please re-formulate!
Pierre: done The log-space you use **log-space** inconsistently sometimes **Log-spaca** and sometimes as **log-space** version of TKF91 reduces the dynamic programming step to finding the greatest of two to three values, depending on the matrix. This in turn meant, that the relative amount of time spent doing so represents a negligible part of the overall algorithm. Nevertheless, we invested some time to devise a vectorized wave-front version of the program.

Initially, we attempted to find a suitable way for storing and accessing the matrix cells, to store data contiguously with respect to the wave-front data access pattern. We concluded that merely indexing **accessing or indexing? Pierre: I**

guess indexing is more specific, to me accessing encompassed indexing. changed it a cell, using the usual row-, and column coordinates, requires a significant amount of computation by itself. Most importantly, indexing a cell what? Pierre: changed it requires significantly more arithmetic operations please quantify! Pierre: we never found out, we abandoned trying to find a indexing formula when we realized that just accessing would require on the order of 20-30 operations, as opposed to the normal 2-3 than the actual computations conducted for each matrix entry.

Furthermore, our preliminary benchmarks had yielded absolutely no performance increase for the vectorized versions of either the basic log-space, or the enhanced, Array-of-Structs implementation what are those major implementations, maybe you should find better names for them Pierre: tried to name them. As a consequence, we abandoned the wave-front vectorization approach.

5 Evaluation and Testing

5.1 Benchmarks and Test Suites

An automated test suite is extremely valuable for re-factoring and optimizing the code, in that it provides the confidence of being able to detect regressions early and get helpful hints on how to fix them. We used the *ctest* module of *CMake* to invoke a Python script driving our console wrapper through a number of inputs. As we also used a continuous integration service running the tests on new commits, we could easily spot when a change accidentally broke the build. This development process led to a high level of satisfaction and faith to try out different ideas in separate branches, while immediately getting feedback.

TODO consistency

When the most obvious optimizations were done and we began to make use of vectorization instructions, we realized a benchmark suite based on the alignment data sets from the University of Oldenburg³ using the Celero C++ benchmarking library⁴ to have a reliable way of validating performance of our different approaches.

5.2 Results of Team I

From each of the six data sets, we picked ten nucleotide strings and computed all 45 pairwise alignments. We found that the median of five *samples*, each of which consisting of the average of ten such *iterations* (i.e. $5 * 10 * 45 = 2250$ alignments per data set) gave reliable performance estimates. As stated in Section 3.5, we performed various performance comparisons based on this benchmark suite, such as the viability of different vector load and store schemes or the SoA, AoS, and AoSoA issue, of which the benchmark results on the reference hardware are depicted in Table 5.2. As already noted in Section 3.5, we deemed the SoA approach the most performant.

³ <http://goo.gl/nlD4nb>

⁴ <https://github.com/DigitalInBlue/Celero>

ms per iteration (compared to best)	SoA	AoS	AoSoA
BDNF	8.30 (1.01)	8.22 (1.00)	8.51 (1.04)
cytb	26.19 (1.00)	29.07 (1.11)	29.65 (1.13)
RAG1	20.88 (1.02)	20.56 (1.00)	20.94 (1.02)
RAG2	24.64 (1.00)	25.09 (1.02)	25.83 (1.05)
RBP3	35.26 (1.00)	36.35 (1.03)	37.62 (1.07)
vWF	44.54 (1.01)	44.10 (1.00)	46.27 (1.05)

Table 1. Benchmark results for SoA, AoS and AoSoA (from left to right)

5.3 Results of Team II

TODO

We conducted all measurements on a desktop machine with 16 GB RAM and a IntelTM i7-2600 CPU with 4 physical cores and hyper-threading.

5.4 Different Logarithm Libraries

Additionally to using the `log` function from the standard C++ header `<math.h>`, we also used the `crlibm` library by Daramy et al. [?]. It includes versions of the natural logarithm computation, which allows the type of rounding to be explicitly specified.

In order to check the correctness of our resulting alignments, we computed the edit distance between the alignments we obtained by using different logarithm implementations. We assigned a cost of 1 to insertions, deletions and substitutions of single letters. We picked all pairs of sequences from the `BDNF_unaligned_sequences.fas` data set from http://www.uni-oldenburg.de/fileadmin/user_upload/biologie/ag/systematik/download/Programs/benchMark_data.tar.gz.

For the remaining parameters, we used $\lambda = 1$, $\mu = 2$, $\pi = (0.27, 0.24, 0.26, 0.23)$ and $\tau = 0.1$. The edit distances were computed in relation to using the logarithm function from `Boost.Multiprecision`.

In particular the `log_ru` function, which explicitly rounds up, produces a notable change in alignment compared to the logarithm function from `Boost.Multiprecision` as can be seen in Table 2. The computed likelihood scores were highly similar and the differences in the alignments were minimal as can be seen in Figure 13. We finally decided to use the `log_ru` function from the `crlibm` library because it had the most consent with the alignments returned by the reference implementation (http://sco.h-its.org/exelixis/web/teaching/practical15/scaledCode/tkf91_scaling.tar.gz).

5.5 Runtime Measurements

The input sequences used to measure execution time were four pairs of randomly generated sequences of length 10, 100, 1000 and 10000 nucleotides. Sequences did not contain ambiguous characters.

Logarithm Function	Average Edit Distance
log from <code><math.h></code>	2.15
log_ru from <code>crlibm</code>	33.15

Table 2. Average edit distances from the alignment returned by using `Boost.Multiprecision` of different logarithm implementations.

Alignment using standard C++ header `<math.h>`

ACGACTAGTCA-GC-TACG-AT-CGA-CT-C-ATTCAACTGACTGACA-TCGACTTA
A-GAG-AGTAATGCATACGCATGC-ATCTGCTATT---CTG-CTG-CAGTGG--T-A

Alignment using `Boost.Multiprecision`

ACGACTAGTCA-GC-TACG-AT-CGA-CT-C-ATTCAACTGACTGACA-TCGACTTA
A-GAG-AGTAATGCATACGCATGC-ATCTGCTATT---CTG-CTG-CAGTGG--T-A

Alignment using `log_ru`

ACGACTAGTCA-GC-TACG-AT-CGA-CT-C-ATTCAACTGACTGACA-TCGACTTA
A-GAG-AGTAATGCATACGCATGC-ATCTGCTATTCTG-CTG-CAGTGG--T-A

Alignment from reference implementation

ACGACTAGTCA-GC-TACG-AT-CGA-CT-C-ATTCAACTGACTGACA-TCGACTTA
A-GAG-AGTAATGCATACGCATGC-ATCTGCTATTCTG-CTG-CAGTGG--T-A

Fig. 13. Difference in alignments using the input parameters $\pi = (0.25, 0.25, 0.25, 0.25)$, $\lambda = 1$, $\mu = 2$, $\tau = 0.1$

We measured the execution times around the kernel of the program, that is excluding any I/O required to load parameters and input sequences, or output of the algorithm. Kernel execution was performed multiple times per input and subsequently averaged. The number of runs executed was dependent on the size of the input, so as to balance accuracy and overall benchmark time (see Table 3).

Length of Sequence	Number of Runs
10	10000
100	1000
1000	100
10000	10

Table 3. Number of runs per sequence length

6 Teaching Results

6.1 What did Team I learn?

While the assigned task was straightforward and very manageable in principle, the implementation details were tricky. Preventing floating point underflow

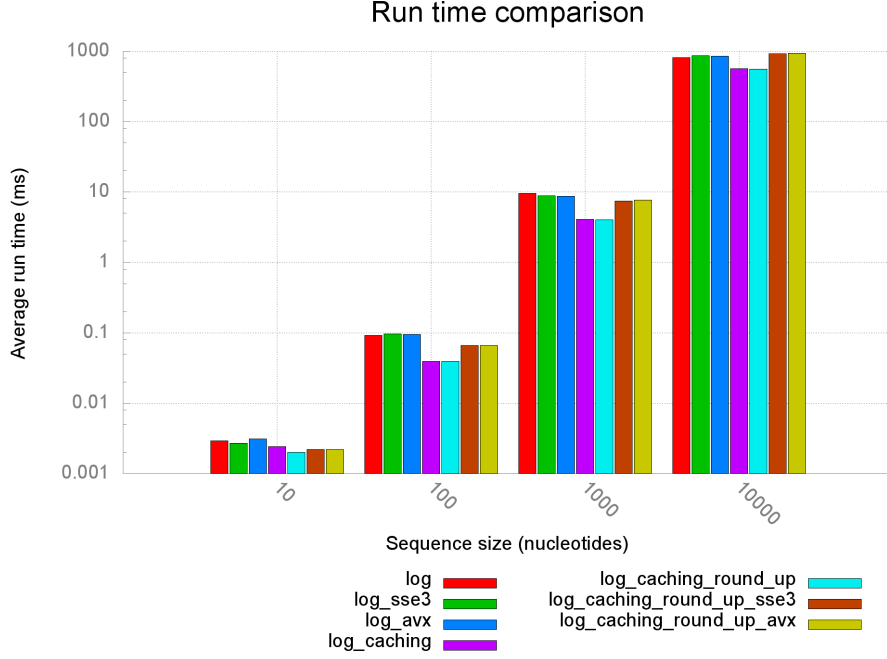


Fig. 14. Comparison of average run times of the different implementations. `sse3` and `avx` post-fixes denote the vectorized versions of the programs, using SSE3 and AVX vector intrinsics respectively (Section 4.3). `log` denotes, that a version built on the log-space transformed version of TKF91 (Section 4.1). Versions using the Array-of-Structs method of storing matrices (Section 4.2) are denoted by the `caching`-keyword, and the `round_up`-keyword denotes the use of the `crlibm` library logarithm function `log_ru` (Section 5.4).

while optimizing performance at the same time represented a challenge. This is because, such technical problems are frequently ignored in other programming practicals at KIT that emphasize functionality. Furthermore, vectorizing with SIMD intrinsics required some rethinking and re-engineering because we needed to identify a suitable memory layout, optimize data alignment, and deal with boundary conditions (padding). While it was satisfying to address and solve problems progressively, there were always additional ideas to further improve the code (w.r.t. performance and design), which made prioritizing tasks important. We enjoyed the satisfaction of gaining yet another percent of execution time in combination with working with SIMD instructions that were not familiar to us.

The basic problem was clearly and outlined. There was enough freedom and time left to work on improving our solution. The scheduled project milestones turned procrastination and last minute work into a non-issue. Most surprisingly, we were always on schedule.

TODO: apart from SIMD were there any other new techniques or tools (e.g. cache analysis, clang compiler etc.) you used for the first time

Well, we didn't do cache analysis and clang wasn't new to us either. We occasionally ran a profiler over our code, but I think that is covered with improving the code w.r.t. performance

6.2 What did Team II learn?

For this assignment we had to tackle two major problems: preventing the loss of precision generated by numerical underflow and pursuing the elusive “most efficient” implementation. In the early stages, we came up with the idea of transforming the algorithm into log-space to solve the numerical issues. This did not only prove to represent an efficient solution, but also allowed us to further simplify the formulas and to pre-compute the vast majority of terms. In this context we also discovered deviations in the output alignments that depend on the specific implementation of the logarithm function being used. We therefore used the `crlibm` library of mathematical functions, that allowed us to assess the impact of logarithm functions with distinct rounding strategies on the final result.

To optimize the code we also used SSE3 and AVX intrinsics to vectorize the most work-intensive portions of the code. While we were unable to produce a faster vectorized code, we did learn how to use vector intrinsics and how to deal with the associated pitfalls. The largest performance gain was attained via a more cache-efficient data structure. For this, we used profiling tools such as `perf` and `valgrind cachegrind` for the first time **correct?**. Finally, we concluded that C++ is appropriate for HPC projects, provided that, the problem and data-structures at hand are well understood.

TODO: apart from SIMD were there any other new techniques or tools (e.g. cache analysis, clang compiler etc.) you used for the first time

6.3 What did the Teacher learn?

TODO

7 Conclusion

For other courses, integration, even faster codes, some of the techniques to be used for more complex stat aligners such as

References

1. Jeffrey Thorne, Hirohisa Kishino, and Joseph Felsenstein. An evolutionary model for maximum likelihood alignment of DNA sequences. *Journal of Molecular Evolution*, 33(2):114–124, 1991.