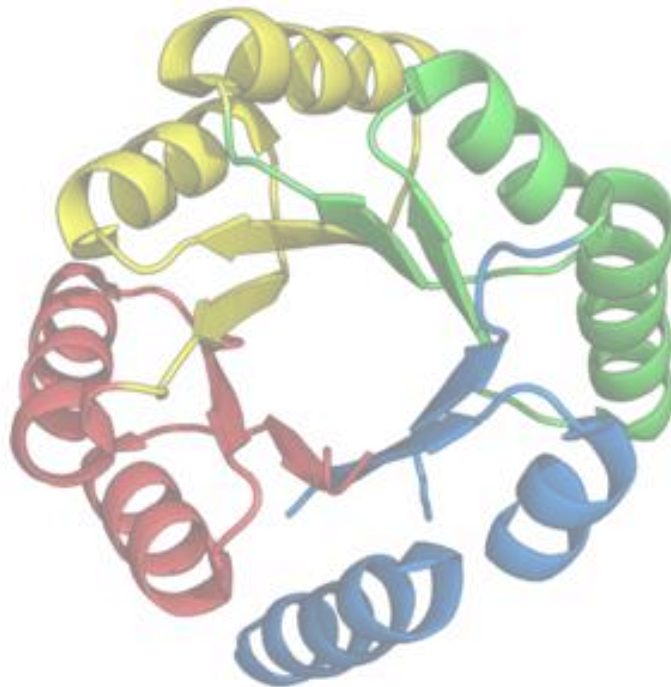**MSc in Data Science and Information Technologies**

National and Kapodistrian University of Athens

# Algorithms in Structural Bioinformatics

Course Project: TIM Barrels

Evangelia Giannaki (7115152100025)
Katerina Mastoraki (7115152100031)
Sophia Nteli (7115152100032)

Athens, 2022

# Contents List

# Introduction

The TIM-barrel fold, also known as alpha/beta barrel, is a conserved protein fold (its sequence is maintained through time and natural selection, despite all the different forces of recombination) and it was named after triose phosphate isomerase, which was the first enzyme that was crystalized ("TIM barrel", Wikipedia, para. 1). The most probable theory about their creation is that they were evolved from a common ancestor through gene duplication and domain fusion. The theory has been accepted by the scientific community as many TIM-barrel proteins possess 2-fold, 4-fold, or 8-fold internal symmetry (Sergio Romero-Romero et al. , 2021).
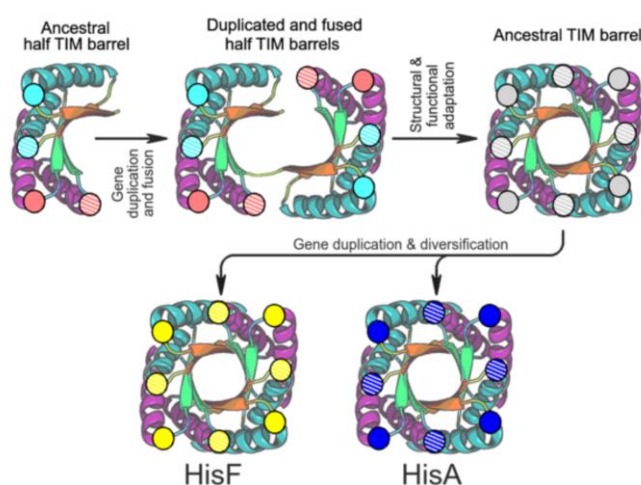


Figure 1: Model for the evolution of TIM barrels through gene duplication and domain fusion

This fold is adopted by approximately 10% of the enzymes and most of them consist of 8 alpha helices and 8 parallel beta strands that alternate along the peptide backbone, connecting between them through loops (Po-Ssu Huang et al., 2015). However, some TIM-barrels contain 7 or even 9 strands and helices. According to CATH classification, each protein contains from 150 to 500 residues, with an average of 298 amino acid residues per structure. Despite the low sequence homology, TIM-barrel structure is well conserved (Nagano et al., 2002). Approximately 160 residues are structurally equivalent among the different proteins while, the remaining residues are located into loop regions. The residues on the loop regions are mainly responsible for the structural maintenance and the enzymatic catalysis of the protein. Thus, each enzyme contains different subsets of residues according to its action and its exact structure and the loop regions are if needed, to contain all these distinct protein domains.

There are a lot of studies about the function and the comparison of the different TIM-barrel families, however none of them focuses on identifying the regions of the proteins which are conserved among them. In the absence of illuminating studies, we were motivated to study TIM-barrels and construct a pipeline for calculating an average structure using numerous TIM barrel PDB files.

# Scope of the project

The aim of this project was the construction of a pipeline to extract a mean TIM-barrel structure.

# Steps of the analysis

## Data collection

The first step of the analysis was the collection of PDB IDs of TIM-barrels. For this purpose, we searched for "TIM-barrels" in the Protein Data Bank (PDB), which returned 8656 PDB entries. To ensure that our dataset contains only proteins with the TIM-barrel fold, we did the same search in CATH database, which returned 1716 results. However, some of CATH's results were not included in PDB's. We decided to continue our analysis only with the common results of both databases, which were 1144 proteins. We used PDB's Advanced Search to select the entries with refinement resolution < 2.5 Å. Then, we sorted them by the number of modeled residues of chain A and we used the first 30 results to proceed with our analysis. The sorting step was crucial to ensure that the proteins used for the analysis will contain only the TIM-barrel fold and not any other structures. The following analysis was also performed for another set of 30 PDB IDs. For this set, before sorting our results, we applied one more filter, which was the "Isomerase" classification of the proteins. We will refer to the two datasets as "IC" (Isomerase Classification) and "NIC" (No Isomerase Classification).
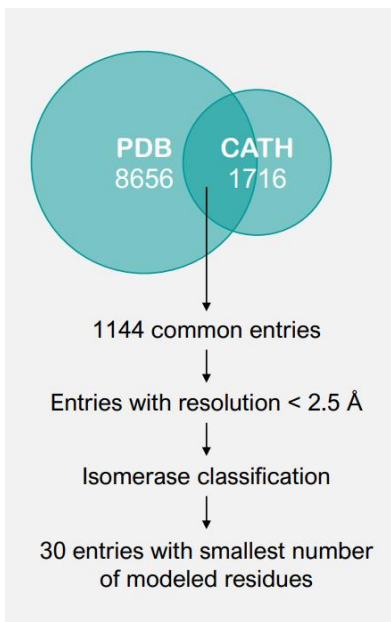


Figure 2: Data collection procedure

## Multiple Structure Alignment

The next step of our analysis was the performance of Multiple Structure Alignment on the two datasets. For this purpose, we used PDBeFold.

### PDBeFold

PDBeFold is a very powerful structure alignment tool which can perform both pairwise and multiple three-dimensional alignment. The structure alignment of the proteins is based on a graph-theoretical approach where the objects in question are represented in graphs, then the graphs are being matched and finally the common subgraphs found are evaluated to form conclusions about similarity.

The algorithm of PDBeFold consists of the following steps (E. Krissinel et al., 2004):

1. The secondary elements (SSEs) of each protein (helices or strands) are represented in a three-dimensional graph
2. The three-dimensional graphs representing the proteins are being matched
3. The individual residues of matched and non-matched SSEs are being mapped by their **Ca** atoms
4. Quality filters to improve the alignment (ex. unmapping Ca atoms of less similar parts) are applied.
5. The quality of the alignment using the Q-score is measured. This function offers a compromise between contradicting requirements of achieving a lower RMSD and a higher number of aligned residues and, therefore, Q is expected to be a more objective indicator of quality of alignment than RMSD and Nalign alone.

$$Q = \frac{N_{align}^2}{\left(1 + \left[\frac{RMSD}{R_0}\right]^2\right) N_1 N_2}$$

To perform the Multiple Structure Alignment, we decided to use only chain A of each protein. Until now, we have the PDB ID's saved in a csv file. However, PDBeFold requires a specific input file format. We have created a Python script that takes as input a csv file with PDB ID's and outputs a text file of the appropriate PDBeFold format.
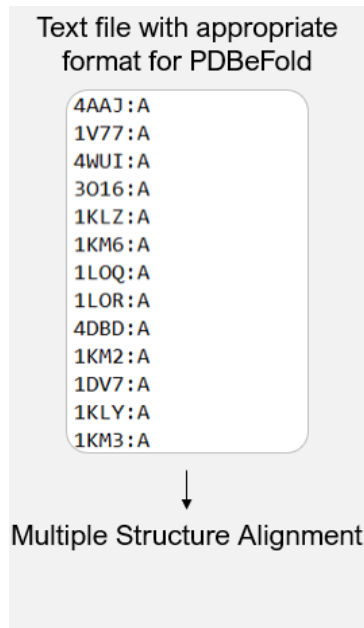
Figure 3: Format of PDBeFold's input

To perform Multiple Structure Alignment, we launch PDBeFold, we select the "multiple" submission form and "List of PDB codes" as source. We then select the file that we generated and submit our query.



Figure 4: Multiple Structure Alignment using PDBeFold

Figures 5 and 6 show the results of the Multiple Structure Alignment for the two datasets. The Q-scores for both our experiments are too low. This means that the quality of the alignment is low, either because the RMSDs are high, or because the number of aligned residues is low. However, the IC's dataset Q-score is extremely low because the filters that we applied were too many and the proteins that we ended up with did not have similar sequence lengths.

For our analysis, we needed more proteins with small sequence length to include only the TIM-barrels fold. In the NIC dataset, we included more proteins with small sequence length, so the final alignment was better because we aligned almost only TIM-barrels.

**Multiple Alignment Results**

| Back to query | Download XML | Download text |

### Summary

| ## | | Structure | $N_{res}$ | $N_{SSE}$ | Consensus scores RMSD | Q-score | | |
|----|---|-----------|-----------|-----------|------|---------|---|---|
| 1 | ☑ | PDB 4aaj:A | 200 | 16 | 2.3336 | 0.2305 | view | download |
| 2 | ☑ | PDB 1v77:A | 202 | 17 | 2.8069 | 0.1953 | view | download |
| 3 | ☑ | PDB 4wui:A | 205 | 16 | 2.4171 | 0.2189 | view | download |
| 4 | ☑ | PDB 3o16:A | 208 | 19 | 1.8577 | 0.2572 | view | download |
| 27 | ☑ | PDB 1u5v:A | 223 | 18 | 2.6925 | 0.1838 | view | download |
| 28 | ☑ | PDB 1z6k:A | 223 | 18 | 2.6742 | 0.1849 | view | download |
| 29 | ☑ | PDB 1ct5:A | 228 | 17 | 2.5776 | 0.1867 | view | download |
| 30 | ☑ | PDB 1b54:A | 230 | 17 | 2.5525 | 0.1866 | view | download |

| Number of aligned residues | 74 | Overall RMSD | 2.362 |
| Number of aligned SSEs | 1 | Overall Q-score | 0.07349 |

Figure 5: Multiple Structure Alignment of the NIC dataset

**Multiple Alignment Results**

| Back to query | Download XML | Download text |

### Summary

| ## | | Structure | $N_{res}$ | $N_{SSE}$ | Consensus scores RMSD | Q-score | | |
|----|---|-----------|-----------|-----------|------|---------|---|---|
| 1 | ☑ | PDB 4aaj:A | 200 | 16 | 1.1361 | 0.0437 | view | download |
| 2 | ☑ | PDB 4wui:A | 205 | 16 | 1.1137 | 0.0429 | view | download |
| 3 | ☑ | PDB 4wd0:A | 230 | 21 | 1.6965 | 0.0329 | view | download |
| 4 | ☑ | PDB 4rcx:A | 235 | 21 | 0.7864 | 0.0398 | view | download |
| 5 | ☑ | PDB 4w9t:A | 235 | 18 | 1.5277 | 0.0338 | view | download |
| 6 | ☑ | PDB 4pc8:A | 239 | 21 | 1.2290 | 0.0358 | view | download |
| 7 | ☑ | PDB 2v2d:A | 241 | 22 | 1.0965 | 0.0366 | view | download |
| 8 | ☑ | PDB 4x9s:A | 241 | 17 | 1.5469 | 0.0328 | view | download |

Figure 6: Multiple Structure Alignment of the IC dataset

## Protein Structural Superposition

The next step of the analysis consists of the batch download of all files that correspond to the PDB IDs that we collected in the previous step by applying several assumptions(NIC or IC protein dataset) to perform structural superposition of the coordinates of each entry.
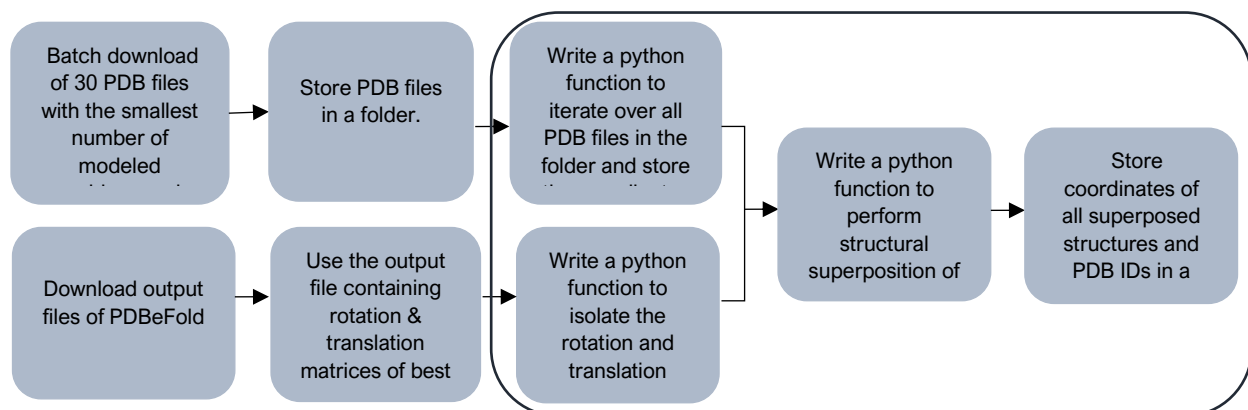


Figure 7: Workflow diagram of protein structural superposition

The batch download was executed in the command line using a shell script that can download multiple PDB archive files by providing a file containing a comma-separated list of PDB IDs. The PDB files were stored in a folder. Then using a Python script, that iterates through each PDB file the coordinates of the Ca atoms of chain A of each entry were isolated. The following step was the usage of the output of PDBeFold that gives information about the rotation-translation matrices of best superposition for each protein of the dataset. The rotation and translation matrices for each protein were isolated using a Python function. To perform structural superposition of the coordinates of each PDB entry we created a Python function that transforms each coordinate by exploiting the information given by the matrices. To be more specific, each coordinate was multiplied by the rotation matrix and then the translation matrix was added to the product.
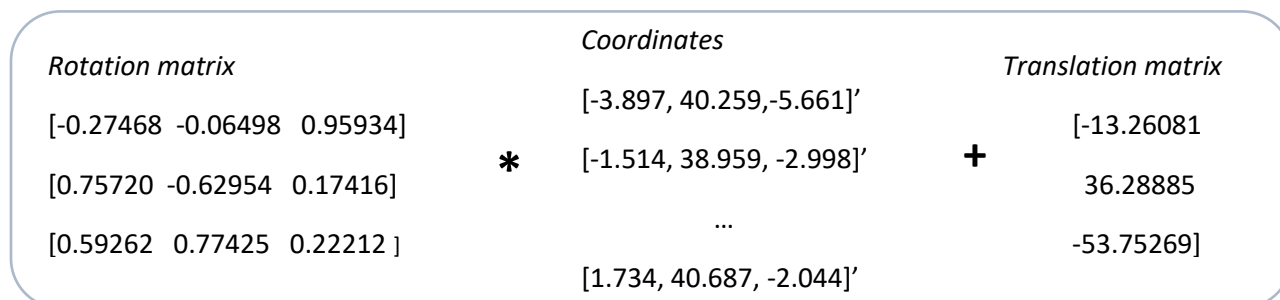


Figure 8: Transformation of the coordinates

The final step of the protein structural superposition was to store the coordinates of all superposed structures in a text file and label them with the corresponding PDB ID, respectively.

## Mean structure calculation

The next step in the pipeline is to calculate the coordinates of the Ca atoms of the mean structure consist of all the TIM-barrel entries of PDB we have used so far. To do that we need:

- The Fasta file of the PDBeFold results, that contains the alignment of every TIM-barrel entry
- The text file with the superposed coordinates for every TIM-barrel entry

The resulted Fasta file of the PDBeFold contains one entry for every PDB file we have used so far in the pipeline. Each entry consists of multiple lines. The first one contains the ID of the PDB file in Protein Data Bank and a small description of the TIM-barrel structure. In the next lines the actual alignment of the structure is depicted. Every position in the alignment may contain:

- A *dash*, a gap in the alignment for the specific residue (position) of the structure, compared with the overall alignment.
- A *lowercase letter*, this residue is aligned with another residue of at least another TIM-barrel structure.
- An *uppercase letter*, this residue is matched with one residue for every TIM-barrel structure.

The text file with the superposed coordinates contains one entry for every TIM-barrel structure. Each entry contains multiple lines. The first line consists of the ID of the PDB file in Protein Data Bank and each of the next lines contains the 3D superposed Ca atoms' coordinates (rotated and translated according to the protein structural superposition mentioned in the previous section) of the structure. Also, before we move on to the calculation of the mean coordinates, we transform the text file with the superposed coordinates, so every number inside a line is separated by the next number with only one space character.



Figure 9: Calculation of the coordinates of the mean structure

To calculate the coordinates of the mean structure we follow two main steps:

- For every entry of the Fasta file we look at a specific position (we start with the first and we move along until no more residues exist). For this specific position we add the superposed coordinates of all the TIM-barrel structures that contains either an uppercase or a lowercase letter.
- When we've checked all the entries for that specific position, we divide with the total number of the PDB entries that contributed to its coordinates' calculation.

Once the previous process is completed, then we save the coordinates into a new text file. This file contains the coordinates for every Ca atom and form the mean structure of all the TIM-barrels we have used in the pipeline. As we have already discussed, we executed the pipeline for two different groups of 30 files each. The text file produced from this step contains:

- 867 3D coordinates – group of isomerase classification (IC protein dataset).
- 446 3D coordinates – group with the smallest residues' number (NIC protein dataset).

## Visualization of the mean structure

The final step of the pipeline is to use the mean coordinates we have calculated to visualize the mean structure and get a 3D representation of the 30 PDB files for the IC and NIC datasets.

For this purpose:

- ✓ UCSF ChimeraX v1.3, a molecular visualization program, was used

- ✓ A PDB file containing the coordinates of the mean structure was created, for the tool to be easier to interpret the results to a 3D structure

To create the PDB file to use it as input to the ChimeraX, we created a Python script that gets as input the text file containing the coordinates of the mean structure. For each line of the text file that corresponds to a triplet of (X, Y, Z) coordinates, we concatenate it with strings that correspond to:

1. The number of the atom, which increases along with the
2. The atom type, which was set to Ca since we have used only carbon atoms of the backbone of the proteins to calculate the mean structure
3. The amino acid type, which randomly was set to SER to all atoms
4. The chain ID, which corresponds to A based on the initial assumptions
5. The residue numbers
6. The coordinates (X, Y, Z) of the mean structure
7. The atom occupancy: most of the time the number is 1.0 (100%) because atoms occupy only one position

8. The atomic displacement parameter, which describes the displacement of atoms from an equilibrium position, and it was randomly set to 37,44
9. The element of the atom, which was set to C

The new lines are written in a text file which simulates a PDB file format, as seen below.



Figure 10: Creation of the PDB file to be used as input to ChimeraX.

Finally, the PDB file of the mean structure is loaded in ChimeraX and as a result we get the 3D visualization of it. Below we observe the representation results of the mean structure extracted from the IC and NIC dataset.
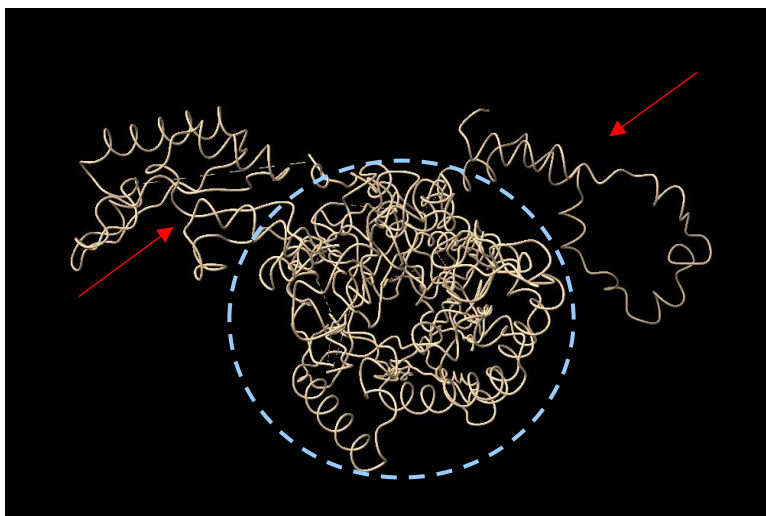


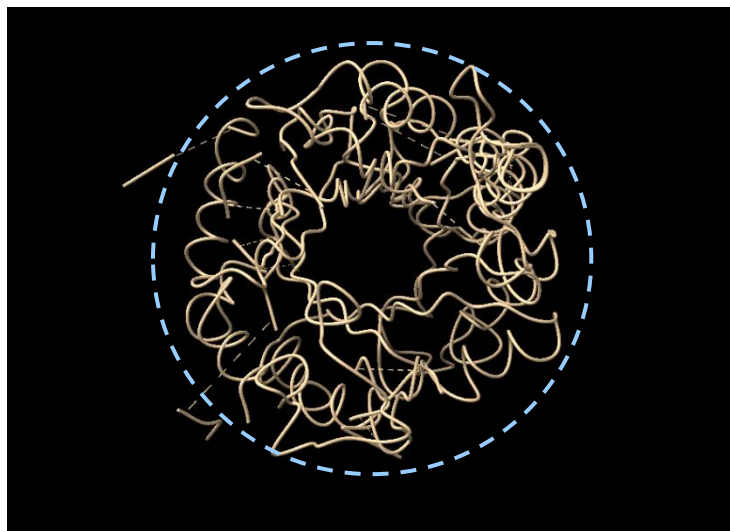Figure 11: 3D visualization of the mean structure of the IC dataset

Figure 12: 3D visualization of the mean structure of the NIC dataset

## Conclusions

Based on the results of our experiments as presented in Figure 11 and 12 we can conclude that the NIC dataset gives a better representation of a mean TIM-Barrel structure than the IC dataset, because it includes almost only structures that belong to the TIM-barrel fold.

To get a better understanding of our results, we launched PDBeFold of the 30 proteins with no isomerase classification along with the mean structure that we have calculated. The goal was to find a minimum and a maximum value of the RMSD score of the structural alignment.

Unfortunately, PDBeFold did not produce any results. Our hypothesis is that either there was an external web - server error in running PDBeFold, or the PDB format that we have created is not compatible with the rest of PDB files.

## Future work

To expand the analysis, we can create a better PDB format of the mean structure and try to re-run the structural alignment using PDBeFold and observe the minimum and the maximum values of the RMSD scores. Furthermore, we can find a distribution of the RMSD values of all pairs of the 30 proteins of the NIC dataset and then check if the median or the mean value of the RMSD values of the mean structure along with the 30 proteins is a good representative.

# Bibliography

1. Nozomi Nagano, Christine A Orengo, Janet M Thornton. (2002). One fold with Many Functions: The Evolutionary Relationships between TIM Barrel Families Based on their Sequences, Structures and Functions. Journal of Molecular Biology, Volume 321, Issue 5.

2. Sergio Romero-Romero, Sina Kordes, Florian Michel, Birte Höcker. (2021) Evolution, folding, and design of TIM barrels and related proteins. Current Opinion in Structural Biology, Volume 68.

3. Kevin T Halloran, Yanming Wang, Karunesh Arora, Srinivas Chakravarthy, Thomas C Irving, Osman Bissel, Charles L Brooks 3rd, C Robert Matthews. (2019). Frustration and folding of a TIM barrel protein. PNAS, Volume 116.

4. 'TIM barrel' (2020). Wikipedia. Available at: https://en.wikipedia.org/wiki/TIM_barrel (Accessed: 22 February 2022).

5. Po-Ssu Huang, Kaspar Feldmeier, Fabio Parmeggiani, D Alejandro Fernandez Velasco, Birte Höcker, David Baker. (2015). De novo design of a four-fold symmetric TIM-barrel protein with atomic-level accuracy. Natural Chemical Biology.

6. E Krissinel, K Henrick. (2004). Secondary-structure matching (SSM), a new tool for fast protein structure alignment in three dimensions. Acta Crystallographica Section D.

7. 'Special relativity' (2018). Wikipedia. Available at: https://en.wikipedia.org/wiki/Special_relativity (Accessed: 5 August 2019).