

THE TIME SLICE SELECTION BAKE-OFF

by

YUYA KAWAKAMI

A THESIS

Presented to the Department of Computer and Information Science
and the Division of Graduate Studies of the University of Oregon
in partial fulfillment of the requirements
for the degree of
Master of Science

June 2022

THESIS ABSTRACT

Yuya Kawakami

Master of Science

Department of Computer and Information Science

June 2022

Title: The Time Slice Selection Bake-off

As the size of data from scientific simulations grows, the ability to identify key time steps in a simulation has emerged as a key challenge. In response, a number of time slice selection methods and algorithms have been proposed. However, no past work has performed a comparative analysis of selection methods as well as their evaluation metrics. This thesis presents results from quantitative and qualitative study of selection methods and evaluation metrics to fill this gap. Our work has three major thrusts. First, we identify similarities and dissimilarities between different time slice selection algorithms. Second, we evaluate conditions under which these methods may fail. Third, we also perform a comparative study with evaluation metrics to investigate how selection methods perform under the set of evaluation metrics in literature. In all, this thesis aims to understand the space of time slice selection methods to inform future research in this area.

CURRICULUM VITAE

NAME OF AUTHOR: Yuya Kawakami

GRADUATE AND UNDERGRADUATE SCHOOLS ATTENDED:

University of Oregon, Eugene, OR, USA
Grinnell College, Grinnell, IA, USA

DEGREES AWARDED:

Master of Science, Computer and Information Science, 2022, University of Oregon
Bachelor of Arts, Mathematics and Computer Science, 2020, Grinnell College

AREAS OF SPECIAL INTEREST:

High-Performance Computing
Scientific Visualization
In Situ Processing

PUBLICATIONS:

Marsaglia, N., Kawakami, Y., Schwartz, S. D., Fields, S., & Childs, H. (2021, October). An Entropy-Based Approach for Identifying User-Preferred Camera Positions. In 2021 IEEE 11th Symposium on Large Data Analysis and Visualization (LDAV) (pp. 73-83). IEEE.

Kawakami, Y., Marsaglia, N., Larsen, M., & Childs, H. (2020). Benchmarking In Situ Triggers Via Reconstruction Error. In ISAV'20 In Situ Infrastructures for Enabling Extreme-Scale Analysis and Visualization (pp. 38-43).

To my Mom, Dad and Sanah.

TABLE OF CONTENTS

Chapter	Page
I. INTRODUCTION	1
II. REVIEW OF TIME SLICE SELECTION METHODS AND EVALUATION METRICS	6
2.1. Domain-agnostic methods	6
2.1.1. Reconstruction-based approaches	6
2.1.2. Information-theoretic approaches	8
2.1.3. Other domain-agnostic approaches	8
2.2. Domain-specific approaches	10
III. EXPERIMENTAL OVERVIEW	12
3.1. Data sets	12
3.2. Selection methods	12
3.2.1. Reconstruction-based approaches	12
3.2.2. Information-theoretic approaches	14
3.2.3. Flow-based* method and uniform selection	15
IV. SELECTION SIMILARITY	17
4.1. Overview	17
4.2. Selection difference	17
4.3. Hierarchical clustering	20
4.3.1. Overview	20
4.3.2. Normalized SD	22
4.3.3. Results	23
4.3.3.1. Results from the first clustering pass (per data set & budget clustering)	23

Chapter	Page
4.3.3.2. Results from the second clustering pass (aggregate clustering)	24
4.3.3.3. Observations on the behavior of the joint entropy selection method	25
V. FAILURE CONDITIONS	27
5.1. Overview	27
5.2. Issues that data sets with high dynamic range encounter	27
5.2.1. Effects on joint entropy selection methods	30
5.2.2. Effects on reconstruction-based selection methods that use Variation of Information	31
VI. METRIC COMPARISONS	35
6.1. Overview	35
6.2. Results	36
VII. CONCLUSION	41
APPENDICES	
A. ADDITIONAL RESULTS FROM EACH SELECTION METHOD	43
B. EVALUATIONS COMPARING TO RANDOM SELECTIONS	49
REFERENCES CITED	52

LIST OF FIGURES

Figure		Page
1.	Variation of Information between two random variables, X and Y , as Venn diagrams.	13
2.	Comparing two selections with Euclidean alignment and Dynamic Time Warping alignment	19
3.	Dendrograms from the first clustering phase	20
4.	Dendrograms from the second clustering phase	21
5.	Volume renderings for the earthquakeMag data set	28
6.	Histograms of data time steps in the earthquakeMag data set	30
7.	Results from each selection method for the earthquakeMag data set at a budget of 10	31
8.	Pairwise mapping cost at heat maps for the earthquakeMag data set	33
9.	Percentile of evaluations among 5000 random selections at a budget of 5.	37
10.	Results from each selection method at a budget of 5 for select data sets. Each dot corresponds to the time step that each selection method makes.	38

LIST OF TABLES

Table	Page
1. Terminology	3
2. Data sets	12
3. Selection methods	16
4. Block sizes for the importance-based approach	16
5. Example of a set of evaluations	35

CHAPTER I

INTRODUCTION

With ever-increasing compute resources at our disposal, the size of scientific data of have also increased, both in the spatial and temporal domains. Coupled with the comparatively slow I/O speeds, it is increasingly difficult to investigate every time step a scientific simulation may produce. As a result, a topic of interest in Scientific Visualization has been time slice selection, i.e., selecting and identifying key time steps from a scientific simulation.

The main motivation of such a system is potential time savings. If done effectively, the selected subset of time steps can be considered as a short summary of the simulation and can remove the need to manually inspect each time step of the simulation after every run. If 10 time steps of a simulation of length 100 suffices to deliver the same or similar insight as the whole time series, then a well-designed algorithm can save (1) the time to run analysis/construct visualizations for the 90 time steps (2) the time to “sort through” the 100 visualizations to extract the main insights.

However, as previous works in this area have highlighted, designing an effective time selection algorithm is difficult. In particular, domain-agnostic selection algorithms, which are the main considerations in this work, are especially difficult as the underlying physics can be vastly different among different fields like astrophysics, climate modeling and fluid dynamics. In contrast, domain-specific algorithms have the advantage that they can exploit domain-specific characteristics of simulations. Work by Bennett et al. [3] is an example of a domain-specific selection algorithm for combustion simulations that tracks and rapid increases in heat releases to identify key time steps in such simulations. However, these domain-

specific methods leverage domain knowledge and are only designed for particular simulation loads.

Another area in which time slice selection is relevant is in situ processing of scientific simulations. In situ processing refers to a processing model that analyzes scientific simulations *as* they run. This contrasts with the more traditional post hoc processing model that analyzes scientific simulations after they complete. Under the post hoc model, a simulation would first dump each time step (or a subset of time steps) to storage. Once the simulation completes, the scientist would inspect the data to extract findings. While post hoc processing is still the dominant model in practice, in situ processing has been growing in interest and popularity over the last decade, primarily due to the I/O constraints on modern supercomputers [13, 7, 6, 16, 17]. To enable in situ processing, inspection routines called “triggers” have been proposed. The main idea behind triggers is to have a lightweight routine that inspects the simulation data that runs *alongside* a simulation. The task for triggers is to decide whether the current time step in a simulation warrants further action as defined by the user: visualization, analysis, further storage, etc.

The main challenge with in situ triggers is that they need to decide whether the current time step warrants further action without the knowledge of future time steps. That is, in situ triggers require time slice selection to be done *as* the simulation runs. Since a well-designed in situ trigger should “trigger” at key time steps of a simulation, one can interpret in situ trigger as a time slice selection algorithm as well - the only difference being whether the selection occurs at the end of the simulation with access to the entire data corpus or not. While the body of past work in in situ triggers is more sparse than post hoc time slice selection algorithms, there are approaches emerging, for example work by Yamaoka et al.

[35]. In all, successful time slice selection for scientific simulation is of significant interest to visualization community. Despite this, no past work has performed a comparative analysis of time slice selection methods, as well as their evaluation metrics. While the lack of comparative analysis of selection methods is troubling for simulation scientists interesting in applying one for their simulation, we believe that the lack of understanding of what each evaluation metric measures is also concerning, especially for future research in this area. Without such an analysis, how should one evaluate their new selection algorithm? Should the community be convinced if the new approach performs well on a metrics A but not on metric B? Even worse, what should one do if two metrics disagree? That is, what should one do if metrics A and B disagree whether one selection is superior to another?

Name	Description
$DATA$	Set of all $data_i$: $\{data_0, data_1, \dots, data_{N-1}\}$
Selection	A subset of $\{0, 1, \dots, N-1\}$, each representing a time step
Budget	The size of the selection. A selection with budget k is a Selection S such that $ S = k$
Selection Method	A function $Method(DATA, k)$ that, given $DATA$ and budget k , makes a corresponding selection S
Evaluation Metric	A function $Metric(DATA, S)$ that, given $DATA$ and selection S , calculates some numeric score. Note that a higher score need not be better.

Table 1. Terminology used in this thesis. This table considers a simulation that runs for N time steps. For any $0 \leq i \leq N - 1$, the data at the i^{th} time step is denoted $data_i$.

Inspired by *The Transfer Function Bake-off* by Pfister et al., [22], with this work, we present the *Time Slice Selection Bake-off*, to explore the space of time slice selection algorithms and evaluation metrics. Before we elaborate on our contributions, we clarify the terminology we employ in this work in Table 1.

In short, the contribution of this work is to illuminate similarities and differences between selection methods and evaluation metrics that are used in practice or proposed in literature. To this end, we perform both qualitative and quantitative analysis of selection methods and evaluation metrics over a range of datasets and budgets. That said, we recognize the space of selection methods, metrics, datasets, and budgets is prohibitively large to explore exhaustively, especially given the large number of selection methods and metrics. One premise of this thesis is that overall trends and characteristics can be understood by considering 10 data sets and 3 budgets. In particular, we use this approach to answer the following questions.

- **RQ1:** Which selection methods produce similar selections?
- **RQ2:** Do we observe notable *outlier* behavior?
- **RQ3:** How do selection methods compare over many metrics?

To address **RQ1**, we perform hierarchical clustering on each selection S , using a distance metric introduced in Chapter IV to identify similar and dissimilar selection methods.

RQ2 will be approached qualitatively by examining the selections that each method makes for each dataset and budget. It is important to note here that, for all of the datasets, we consider the “correct” selection for a particular budget to be unknown. Without the appropriate domain experts, it is difficult to establish an appropriate “ground truth,” even if such a selection exists. (We would argue that visualization experts and domain scientists would reasonably disagree on what an ‘optimal’ selection is for a given dataset and budget.) However, by focusing on the trends in the selections, we believe we can characterize the selection methods as well as point out clear failures conditions.

To address **RQ3**, for every selection method SM , we evaluate its selection S with every selection metric EM to identify if any method SM performs favorably across the set of metrics EM .

CHAPTER II

REVIEW OF TIME SLICE SELECTION METHODS AND EVALUATION METRICS

A number of past works have considered time slice selection methods for scientific simulations. This section reviews the past body of work in this area.

§2.1 reviews domain-agnostic approaches: reconstruction-based approaches (§2.1.1), information-theoretic approaches (§2.1.2) and other domain-agnostic approaches (§2.1.3). Finally, §2.2 reviews domain-specific approaches.

2.1 Domain-agnostic methods

2.1.1 Reconstruction-based approaches. The objective for reconstruction-based approaches is find a selection S , such that it minimizes a particular reconstruction error. The two main components for this set of methods are:

1. *Method for reconstruction:* Given some selection S , the reconstruction method generates reconstructed $data'_i$ for every $data_i \in DATA$ with access only to the data in the time steps in S .
2. *Error function between the original and reconstructed fields:* Given the original field $data_i$ and reconstructed $data'_i$, the error function quantifies the difference between the two.

Given these two components, the reconstruction-based approaches defines the *cost* of a selection S as follows. For some selection S , an error function f and the reconstructed set of $data'_i$:

$$cost = \sum_{0 \leq i \leq N} f(data_i, data'_i) \tag{2.1}$$

Finally, reconstruction-based approaches consider a selection S of budget k *optimal*, if $|S| = k$ and if it minimizes Eq. 2.1.

The most common method for reconstruction is linear interpolation, where the time slices in selection S are used to reconstruct every $data'_i$ by linearly interpolating between the closes available time steps. Zhou and Chiang used linear interpolation as the reconstruction method and proposed Variation of Information (VI) as the error function [37]. VI is a metric derived from information theory and measures the sum of *information* present in X after observing Y and in Y after observing X . In addition to VI, Zhou and Chiang also considers root-mean-square error (RMSE).

Alternatives to linear interpolation have also been proposed. Tong et al. proposed Dynamic Time Warping (DTW) as a method of reconstruction as a nonlinear time mapping method [27]. In this work, they employ Earth Mover’s Distance (EMD) and isosurface similarity map[9] as the error function, but state that these were specifically chosen for the datasets that they considered (i.e., not domain agnostic). However, their DTW-based method does not depend on a particular error function; therefore, it can be used broadly with any error function.

For a given budget, the *optimal* selection for reconstruction-based approaches can be calculated through dynamic programming. That said, the time complexity for these methods grow cubically with respect to the number of time steps which can become impractical for large datasets. Hence, some have also proposed approximation methods to decrease the time complexity [37].

Reconstruction-based approaches are common in evaluation metrics as well. As an evaluation metric, the cost from Eq. 2.1 is directly used for any selection S . The following three works have used a reconstruction-based approach for

evaluation, all using linear interpolation as the method for reconstruction. First, Kawakami et al. proposed an evaluation framework for in situ triggers using linear interpolation along with L1 norm of difference as the error function [11]. Second, Porter et al. evaluated their selection method using linear interpolation along with and root-mean-squared-error (RMSE) and peak signal-to-noise-ratio (PSNE) as the error function [23]. Third, Pulido et al. evaluated their selection method using linear interpolation along with image quality metrics like structural similarity index measure (SSIM) and universal quality index (UQI), among others like RMSE, as the error function [24].

2.1.2 Information-theoretic approaches. The field of scientific visualization has incorporated concepts from information theory many times to solve challenges [30, 5], including time slice selection. Previously mentioned work by Zhou and Chiang that used variation of information as the error function is one example [37]. Wang et al. formalized the notion of *importance* values, based on a sliding time window and conditional entropy values with adjacent time steps. By using applying this metric on a block-partitioned dataset, they proposed to select time slices by maximizing joint entropy of the adjacent time steps in selection S [31]. Shannon’s entropy, which measures the amount of “information” in data, has also been proposed for in situ triggers, where trigger would ‘fire’ if the change in entropy from the last trigger instance exceeds some predefined value [13].

2.1.3 Other domain-agnostic approaches. Frey and Ertl presented a flow-based approach for time slice selection, where the difference between two time steps is quantified via a *flow-based* metric [10]. They quantify the difference between two time slices in three steps: (1) probabilistically sample from data volume, (2) construct a flow graph between the two samples, and (3) calculate the

minimum-cost flow in the said graph. Finally, a selection S is made such that it minimizes the difference with respect to the original time series. Ling et al.[14] presented a time slice selection method that relied on detecting spatial deviations and temporal deviations for simulations running in a multi-process setting. Their work had each processor to maintain a KDE-estimated PDF of the variables in the simulation. They used an ensemble of decision tree regressor to detect whether (1) data in a particular processor exhibited different behavior than the rest (spatial deviations) (2) data at current time step exhibited different behavior than the previous time step (temporal deviations). Pulido et al. developed a time slice selection method based on non-negative Tucker factorization (NNTF) [24]. The key idea for their work was to use a non-negative Tucker decomposition on the 4D tensor (x,y,z,time) and extract “influential” time points as the key time steps. Their use of non-negative Tucker decomposition involved representing the simulation data as a 4D tensor (x,y,z,time) and extracting “influential” time steps as the time slices selection.

Some selection methods have specifically considered an *in situ* setting. Myers et al. presented a time step selection algorithm that fitted a piecewise linear model to the data stream as the simulation ran [19]. Each incoming time step would be evaluated on the linear model, until the precision of the previous model drops to predefined level. Once that happens, the algorithm would start a new linear model, and declare that time step as *salient*. Yamaoka et al. developed an in situ time slice selection algorithm that leveraged KL divergence to detect changes in the state of the simulation [35]. Like Ling et al., they target a multi-process simulation run, and attempt to track changes in the simulation by calculating the

KL divergence between the data in the current time step and the last time step when the algorithm made a selection.

There is also a growing body of work utilizing deep learning to time slice selection [23, 36, 15]. Porter et al. trained an encoder-decoder network to implicitly learn feature descriptors of each time step in a latent space [23]. Once the network is trained, they perform further dimensionality reduction from 1024 to 2 dimensions via t-SNE, before making selections based on the “path” that each time step forms in the projected 2D space. Zhang et al. developed a time slice selection method based on a neural network that, given data from two time slices, reconstructed the data for the intermediate time steps [36]. Based on this data reconstruction module and a predefined error bound, they performed time slice selections *in situ*. However, it should be noted that their method specifically targets ensemble simulations, where the data reconstruction module can be trained on a simulation data with similar characteristics prior to use.

2.2 Domain-specific approaches

In contrast to domain-agnostic approaches, domain-specific time slice selection approaches can incorporate prior knowledge of the simulation and its dynamics. Bennett et al. presented a time slice selection method for combustion simulation by predicting rapid heat releases in these simulations [3]. This work was continued by Salloum et al., who proposed a new metric leveraging the coefficient of variation to capture rapid heat release in combustion simulation that was more accurate and increased the robustness of trigger detection [26]. Banesh et al. presented a method to detect mesoscale ocean eddies in large simulated or observational ocean data [2]. Their method was an extension to the previously introduced work by Myers et al. [19] that utilized a piecewise linear model to model

the simulation. Liu et al. utilized a Siamese network comprised of two CNNs to identify key time steps in flow field data for Computational Fluid Dynamics (CFD) simulations [15]. A key idea in their work to use a Siamese network to assess the similarity between data in two time steps. While their method outperformed the Myers' method [19] for their CFD data, NN-based methods like this method and the ones introduced in §2.1.3 need to be trained before being applied, which can be time consuming especially as the data sizes increases.

CHAPTER III

EXPERIMENTAL OVERVIEW

3.1 Data sets

This work consider 10 different datasets, which are described in Table 2.

Name	Variable of Interest	Description
Asteroid	tev, Temperature in electronvolt (eV)	Simulation of an asteroid impact with ocean $300 \times 300 \times 300$, 237 time steps [21]
Cloud	cct, Cloud top pressure in Pascals	Atmospheric simulation of clouds in Germany $1429 \times 1556 \times 1$, 481 time steps [8]
Droplet	NA, Only one field	Simulation of two droplets colliding $256 \times 256 \times 256$, 167 time steps [4]
EarthquakeMag	Magnitude of velocity	Earthquake simulation $750 \times 375 \times 100$, 227 time steps [1]
HurricanePressure	Pressure in Pascals	Atmospheric simulation of Hurricane Isabel $500 \times 500 \times 100$, 48 time steps [20]
MantleTemperature	Temperature in Kelvin	Earth mantle convection simulation $180 \times 201 \times 360$, 251 time steps [25]
Jet	NA, Only one field	Turbulent combustion simulation $480 \times 720 \times 120$, 121 time steps
Bottle	NA, Only one field	Laser pulse shooting through a bottle 900×430 , 465 time steps [29]
Cloverleaf	Energy	Hydrodynamics simulation $260 \times 132 \times 260$, 91 time steps [28]
Radiation	Gas temperature in Kelvin	Cosmology simulation of radiation waves $600 \times 248 \times 1$, 200 time steps [34]

Table 2. Data sets

3.2 Selection methods

The selection methods we consider in this thesis can be broadly be separated into two categories: (1) Reconstruction-based approaches and (2) Information-theoretic approaches. In addition, §3.2.3 introduces two methods based neither on reconstruction nor information theory.

3.2.1 Reconstruction-based approaches. Reconstruction-based selection methods were introduced in §2.1.1. Recall that there are two components to reconstruction-based selection methods: (1) The method to reconstruct every

$data'_i$ and (2) error function, f , used to evaluate the dissimilarity between $data_i$ and $data'_i$. This thesis considers two reconstruction methods: linear interpolation for every mesh location (LI) and Dynamic Time Warping (DTW)[27]. Furthermore, four error functions are considered: L1 norm of difference (L1N), L2 norm of difference (L2N), mean squared error (MSE) and Variation of Information (VI) using 128 bins as used in [37]. Intuively, Variation of Information measures the difference of *information* contained between two random variables X and Y . The Venn diagram in Fig. 1 provides some graphic intuition into $VI(X, Y)$. A key property for reconstruction-based approaches is that it enforces that the first and final time step be included in the selection. This is to ensure that the reconstruction methods have the sufficient information to reconstruct the set $data'_i$ for every i . Lastly, recall from §2.1.1 that these selections can be made via dynamic programming. In total, $2 \times 4 = 8$ reconstruction-based selection methods are considered.

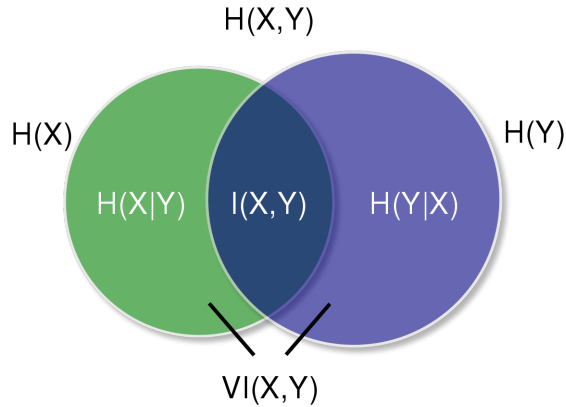


Figure 1. Variation of Information between two random variables, X and Y , as Venn diagrams.

3.2.2 Information-theoretic approaches. The information theoretic approaches we consider are the joint entropy approach and the importance-based selection introduced by Wang et al. [31]

The joint entropy method attempts to maximize the joint entropy in selection S . For some selection $S = \{s_0, s_1, \dots, s_{n-1}\}$, the joint entropy method defines the quality of S of budget n as in Eq. 3.1.

$$\text{Selection quality} = \sum_{0 \leq i \leq n-2} JE(\text{data}_{s_i}, \text{data}_{s_{i+1}}) \quad (3.1)$$

The selection S that maximizes Eq. 3.1 is considered *optimal* under the joint entropy method.

Eq. 3.1 was used as a evaluation metric in [32], but, of course, can also be used as a selection method by maximizing the joint entropy. In a similar fashion to the calculations for reconstruction-based approaches, a dynamic programming method can be employed to calculate the best selection for some budget k .

The importance-based approach by Wang et al. first subdivides the data into equal sized blocks. Next, the *importance* of each block is defined as the sum of conditional entropies of the block with those at the same spatial location in adjacent time steps. Specifically, for some data block X_j at time t , Wang et al. defines its *importance* as seen in Eq. 3.2, where $Y_{j,t}$ refers to the j th block at time step t . Finally, the *importance* of a time step is defined as the sum over all blocks as in Eq. 3.3.

$$A_{X_j,t} = 0.5 \cdot H(X_{j,t}|Y_{j,t-1}) + 0.5 \cdot H(X_{j,t}|Y_{j,t+1}) \quad (3.2)$$

$$A_t = \sum_i A_{i,t} \quad (3.3)$$

With the *importance* of each timestep, Wang et al. then partitions the time series into contiguous sections with approximately equal sum of *importance* values. Finally, a time step from each contiguous partition is chosen such that the joint entropy of every selection in S is maximized. Hence, the importance-based method can be thought of a variant of the joint entropy methods, except that it restricts the pool from which the selections can be made. The size of each block used in this paper for each data set is listed in Table 4 and details of the method can be found in [31].

Both joint entropy and importance-based selection methods requires a histogram to be constructed for each time step, thus requires the number of bins to be defined. In this paper, we consider four bin counts: 8, 64, 128, and 256 bins. Therefore, $2 \times 4 = 8$ total information-theoretic methods are considered.

3.2.3 Flow-based* method and uniform selection. In this work, we also consider the flow-based method by Frey and Ertl [10]. However, we make one key modification: we enforce that the first and last time step be included in the selection. For any budget k , recall that reconstruction-based methods effectively have $k - 2$ degrees of freedom. To ensure that comparisons are valid across methods with the same budget, this paper considers a modified flow-based* method that enforces the inclusion of the endpoints. We also consider the uniform selection which places selections evenly throughout the time series.

Also note that a number of selection methods can also be directly used as an evaluation metric as well. The full set of selection methods this work considers and whether they are also valid metrics are shown in Table 3.

Name	Description	Can be a evaluation metric
<i>LIL1N</i>	Linear interpolation with L1 norm of difference	Yes
<i>LIL2N</i>	Linear interpolation with L2 norm of difference	Yes
<i>LIMSE</i>	Linear interpolation with mean squared error of difference	Yes
<i>LIVI</i>	Linear interpolation with Variation of Information	Yes
<i>DTWL1N</i>	Dynamic time warping with L1 norm of difference	Yes
<i>DTWL2N</i>	Dynamic time warping with L2 norm of difference	Yes
<i>DTWMSE</i>	Dynamic time warping with mean squared error of difference	Yes
<i>DTWVI</i>	Dynamic time warping with Variation of Information	Yes
<i>JE8</i>	Joint entropy with 8 bins	Yes
<i>JE64</i>	Joint entropy with 64 bins	Yes
<i>JE128</i>	Joint entropy with 128 bins	Yes
<i>JE256</i>	Joint entropy with 256 bins	Yes
<i>Importance8</i>	Importance-based selection with 8 bins	No
<i>Importance64</i>	Importance-based selection with 64 bins	No
<i>Importance128</i>	Importance-based selection with 128 bins	No
<i>Importance256</i>	Importance-based selection with 256 bins	No
Flow-based*	Modified version of the flow-based method [10]	Yes
Uniform	An evenly-spaced selection	No

Table 3. Selection methods

Data set	Block size
Asteroid	$20 \times 20 \times 60$
Cloud	$50 \times 50 \times 1$
Droplet	$16 \times 16 \times 16$
EarthquakeMag	$50 \times 30 \times 10$
HurricanePressure	$25 \times 25 \times 10$
MantleTemperature	$20 \times 20 \times 20$
Jet	$24 \times 26 \times 20$
Bottle	$45 \times 25 \times 1$
Cloverleaf	$23 \times 20 \times 23$
Radiation	20×20

Table 4. Block sized used for the importance-based approach [31]

CHAPTER IV

SELECTION SIMILARITY

4.1 Overview

In short, this section addresses **RQ1**: which selection methods produce similar selections? In other words, we want to cluster the selection methods in Table 3 into those that exhibit similar behavior. The importance of such an analysis is to gain a better understanding of how selection methods relate to each other across the range of data sets and budgets. To this end, we apply two passes of agglomerative clustering.

1. Agglomerative clustering on the result that each selection method makes for a particular data set *and* a particular budget.
2. Agglomerative clustering on the result that each selection method makes across *all* data sets for a particular budget.

For example, the first clustering pass would cluster selection methods based on its results on the `asteroid` data set for a budget of 5, while the second pass would cluster selection methods on its results on *all* data sets for a budget of 5. In order to quantitatively cluster selections using a hierarchical clustering model, a metric is required to capture the dissimilarity between two selections.

4.2 Selection difference

Given two different selections of budget k , S_1 and S_2 , we want to define a *selection difference* $SD(S_1, S_2)$, to capture the dissimilarity between the two. In order to define SD , we borrow ideas from alignment-based metrics based on time series analysis. In particular, consider two alignment schemes: Euclidean alignment and Dynamic Time Warping alignment. (An important note: The Dynamic Time Warping (DTW) *alignment* introduced here is entirely different than the DTW-

based selection method by Tong et al [27]. Though they are both based on the same idea from time series analysis, the two should not be confused.)

For two time slice selections of budget k , $S_1 : \{x_0, \dots, x_{k-1}\}$ and $S_2 : \{y_0, \dots, y_{k-1}\}$, an *alignment* between S_1 and S_2 is a set $P = \{p_0, \dots, p_m\}$ of pairings p_i . Each p_i is pair of elements (x_i, y_j) : one selection from S_1 , and another from S_2 .

A simple metric we can devise is to consider the difference between S_1 and S_2 as the summation of pairwise difference at each index. Aligning every i^{th} element with each other between two time series is referred to as the Euclidean alignment[33] (i.e. for any i , $p_i = (x_i, y_i)$). The difference metric that arises from this alignment can be written as follows:

$$\text{Euclidean alignment cost}(S_1, S_2) = \sum_i \text{abs}(S_1[i] - S_2[i]) \quad (4.1)$$

However, as noted by Keogh and Pazzani, the Euclidean alignment between time series is highly sensitive to small distortions in the time axis [12], thus, making it less appealing in many scenarios.

Thus, the more commonly used alternative, and the method used in this paper, is the Dynamic Time Warping (DTW) alignment [12], which allows for a more flexible and intuitive dissimilarity metric between two time series. A valid DTW alignment satisfies the following three conditions.

1. The first and last indices are paired together. In other words, $(x_0, y_0) \in P$ and $(x_{k-1}, y_{k-1}) \in P$.
2. Every x_i, y_j appears at least once in a pairing in P .
3. The pairings between S_1 and S_2 must be monotonically increasing. In other words, for any $x_i, x_j \in S_1$ with $i < j$, its respective pairing $y_s, y_t \in S_2$ must obey $s \leq t$, and vice versa. The monotonicity of pairings is enforced to preserve the time ordering in the pairings.

Finally, the optimal DTW alignment is the set of pairings P that minimizes the difference in each pairing p_i . Formally, for two selection S_1, S_2 , and the set of all valid DTW alignments, denoted \mathcal{P} , the DTW alignment cost is defined as follows.

$$\text{DTW alignment cost}(S_1, S_2) = \min_{P \in \mathcal{P}} \sum_i \text{abs}(P_i[0] - P_i[1]) \quad (4.2)$$

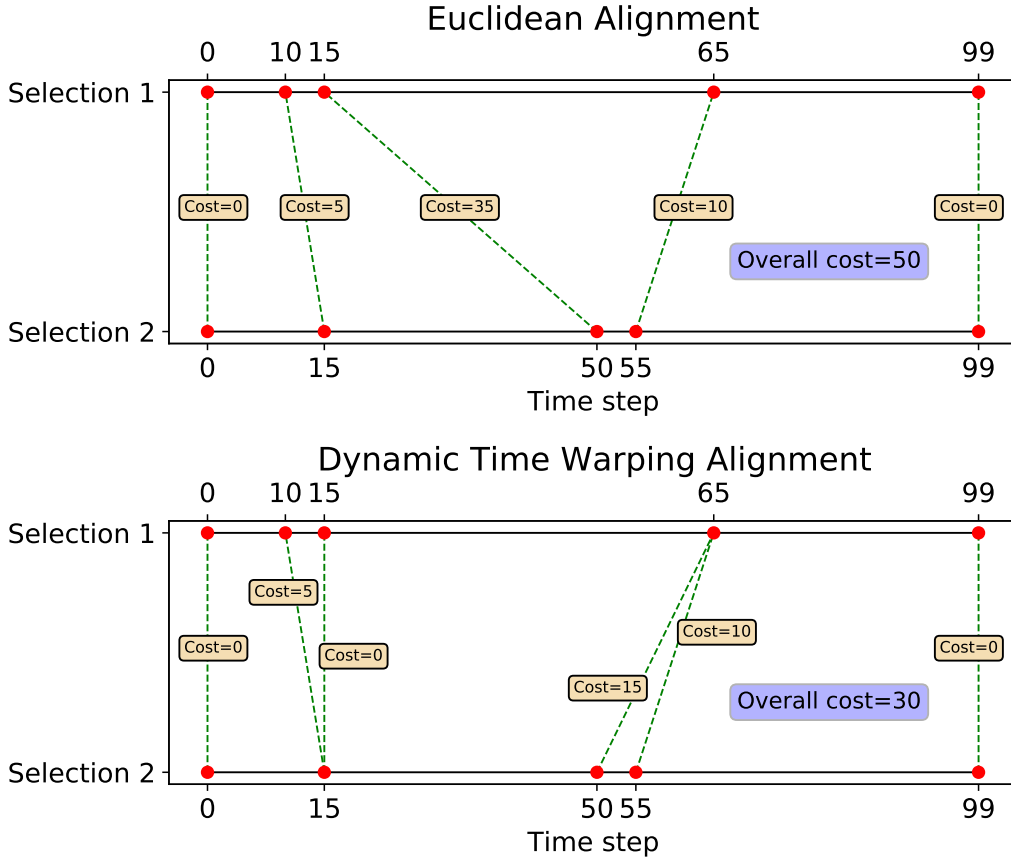


Figure 2. Euclidean alignment (top) and Dynamic Time Warping (bottom) alignment between two selections: $\{0, 10, 15, 65, 99\}$ and $\{0, 15, 50, 55, 99\}$. The dotted lines denote the pairings that would be produced under the respective alignment schemes. Notice that Euclidean alignment has cost of $0+5+35+10=50$, while the DTW alignment has cost of $0+5+0+15+10+0=30$.

Fig. 2 plots a notational relation between the Euclidean alignment cost and DTW alignment cost. As Keogh and Pazzani pointed out, the Dynamic

Time Warping alignment better describes the difference between two time slice selections. Here, the Euclidean alignment that forces the (15, 50) pairing paints a distorted picture. As a result, as mentioned prior, in this thesis, we employ the DTW alignment cost as $SD(S_1, S_2)$.

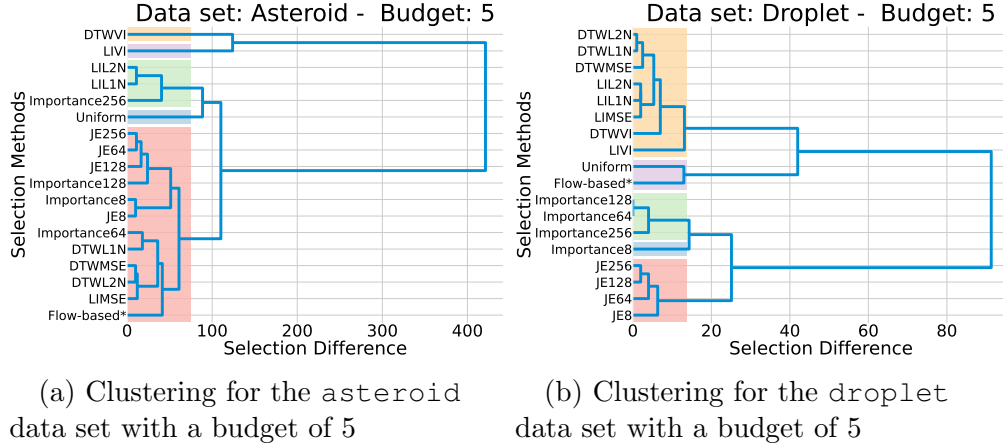


Figure 3. Dendrograms from the per data set/budget hierarchical clustering of selection methods (i.e. the first clustering phase). Each clustering is also color-coded based on results if 5 clusters were to be extracted.

To calculate $SD(S_1, S_2)$, dynamic programming can be used, following the recurrence relation in Eq. 4.3. Let $DP[i][j]$ denote the *selection difference* between the first i elements of S_1 and first j elements of S_2 . Then, the following holds.

$$DP[i][j] = \min(\begin{aligned} &abs(S_1[i] - S_2[j]) + \\ &DP[i-1][j], DP[i-1][j-1], DP[i][j-1] \end{aligned}) \quad (4.3)$$

Finally, notice that $DP[k][k]$ is equal to $SD(S_1, S_2)$ for any two selections S_1, S_2 of budget k .

4.3 Hierarchical clustering

4.3.1 Overview. For a particular data set and budget, we determine similarities across the selection methods by applying agglomerative clustering with this SD metric. (Note that average-linkage clustering is used as the method to join

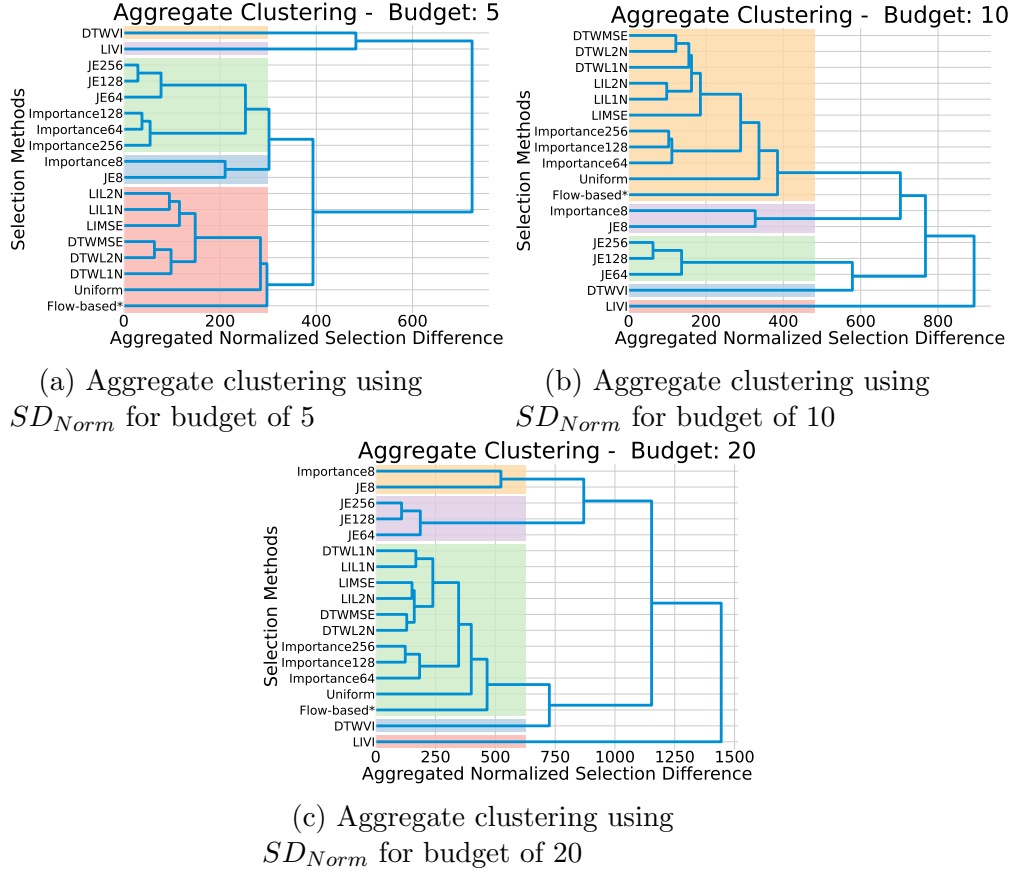


Figure 4. Dendrograms from the aggregate hierarchical clustering of selection methods (i.e. the second clustering phase). Each clustering is also color-coded based on results if 5 clusters were to be extracted.

clusters in the hierarchical clustering process as recommended by Manning et al. [18].) However, this per data set/budget clustering only highlights local clusterings. Therefore, as noted in §4.1, in order to cluster the behavior of selection methods across data sets for a particular budget, we use agglomerative clustering on the aggregated selection differences across all data sets.

There are number of ways to aggregate selection differences across data sets. The most simple method, and perhaps obvious, is a simple addition of SD values. In other words, we can assess the difference between two selection methods for some budget k by summing SD value for each of the 10 data sets. However, simple

addition of SD across data sets can be a problem, since the number of time steps in each data set varies. The range of values that SD can take on increases for data sets with longer time series; hence, simple addition of raw SD values will inflate contributions from data sets with more time steps. Therefore, for the aggregate hierarchical clustering phase, we devise a normalized SD metric, to overcome this issue.

4.3.2 Normalized SD. In essence, for the normalized SD , or SD_{Norm} , we consider each element in a selection, not as a time slice index, but as a percentage of the whole time series. By way of example, consider a data set with 50 time steps, and a selection for this data set $\{0, 10, 15, 49\}$. In the SD_{Norm} calculation, we instead treat this selection as $\{0\%, 20.4\%, 30.6\%, 100\%\}$ or $\{0, 20.4, 30.6, 100\}$, and proceed with the same calculation as in Eq. 4.2. This transformation is permissible for a several reasons. First, notice that this linear transformation does not affect the local per data set/budget clustering from the first clustering phase, since SD and SD_{Norm} would bring rise to the same clusters in the first per data set/budget clustering phase. Second, the number of time steps in a particular data set is often arbitrary. Simulation scientists will often decide the rate to dump the simulation data to disk (e.g. every 5 simulation cycles). Therefore, we argue the actual time step index of the selection is less interesting than where it is temporally located in the time series.

The key, desirable property of SD_{Norm} is that it sidesteps the issue that SD had under simple addition across data sets. Since the values are all within $[0, 100]$, simple addition of SD_{Norm} can be used to cluster the behavior of selection methods across all data sets. (We cannot, however, use SD_{Norm} to cluster across different budgets, since the range of SD_{Norm} increases as the budget increases.)

With SD_{Norm} established, we define the pairwise difference between two selection methods SM_1, SM_2 for a budget k in the aggregate clustering phase in Eq. 4.4. Let D denote the set of all considered data sets. Then,

$$\text{Aggregated difference}(SM_1, SM_2, k) = \sum_{dset \in D} SD_{Norm}(SM_1(dset, k), SM_2(dset, k)) \quad (4.4)$$

4.3.3 Results.

4.3.3.1 Results from the first clustering pass (per data set & budget clustering). For this selection similarity analysis, we consider three budgets: 5, 10, and 20. This analysis does not consider larger budgets for two reasons: (1) As described in the introduction, time slice selection algorithm are often used to generate a short summary of the data set. Therefore, comparisons at lower budgets are more meaningful. (2) Differences between selection methods are more accentuated at lower budgets. Now, since there are 10 data set under consideration, this yields $10 \times 3 = 30$ per data set/budget clusterings, and 3 aggregate clustering using SD_{Norm} for each budget. Instead of analyzing all 33 clustering results and dendrograms, we highlight five clustering results and show their dendrograms in Figs. 3 and 4.

Figs. 3a and 3b consider the asteroid and droplet data set, respectively, at a budget of 5. Notice that the clustering results between the two share some similarities, but also have key differences. Some observations follows.

- The joint entropy methods (i.e. $JE8$, $JE64$, $JE128$, and $JE256$) make similar selections in both data sets.
- DTW-based methods except $DTWVI$ (i.e. $DTWL1N$, $DTWL2N$, $DTWMSE$) make similar selections for both data sets.

- The relative behavior of the importance-based method (i.e. *Importance8*, *Importance64*, *Importance128*, *Importance256*) is different. While their behavior are similar for the `droplet` data set, the behavior is more varied for the `asteroid` data set.
- LI-based methods except *LIVI* (i.e. *LIL1N*, *LIL2N*, *LIMSE*) make similar selections for the `droplet` data set, but for the `asteroid` data set, *LIMSE* performs differently than the other two.

4.3.3.2 Results from the second clustering pass (aggregate clustering). With the 30 local clusterings as mentioned prior, however, exploring each clustering and their differences is (1) time consuming and (2) less interesting than characterizing general behavior over all datasets. Instead, this section will mainly focus on the results from Fig. 4 to answer the question posed in §4.1 - which selection methods produce similar results? In the aggregate, we make the following observations.

- *LIL1N*, *LIL2N*, *LIMSE*, *DTWL1N*, *DTWL2N*, and *DTWMSE* all produce similar results for all data sets for all budgets.
- *JE64*, *JE128*, *JE256* produce similar results for all data sets for all budgets.
- *Importance64*, *Importance128*, *Importance256* produce similar results for all data sets for all budgets.
- Flow-based* generally performs most similar to *LIL1N*, *LIL2N*, *LIMSE*, *DTWL1N*, *DTWL2N*, and *DTWMSE*.
- *JE8* and *Importance8* are most similar to each other, and generally dissimilar to the rest.
- VI-based methods (i.e. *LIVI*, and *DTWVI*) are outliers. *LIVI* forms a single leaf at budgets 10 and 20. At budget 5, though it forms a cluster with

DTWVI, LIVI is still significantly different that the rest of the selection methods. Despite it following the same reconstruction-based approach as methods like *LILIN*, their behaviors present a significant departure from them. The reason for this departure is explored in Chapter V.

4.3.3.3 Observations on the behavior of the joint entropy

selection method. A key difference that separates reconstruction-based and joint entropy selection method is whether they enforce some level of temporal spacing in their selections. Reconstruction-based selection methods, by way of their design, tend to prefer selections that are temporally spaced out. In fact, this naturally follows from how the cost for reconstruction-based methods are defined. Regardless of the specific reconstruction paradigm, these methods hinge on finding time steps in the simulation from which reconstruction of other time steps are best achieved. Therefore, skipping many intermediate time steps tends to introduce high error in the cost calculation, leading it to prefer to spread out its selections. In contrast, joint entropy methods pay little attention to how selections are “spread” out temporally, since no reconstruction is done. Thus, there is no inherent penalty for skipping many intermediate time steps in adjacent selections - they only serve to maximize the joint entropy of adjacent selections. This is seen in Fig. 7 and Figs. 10b and 10c. Notice that the joint entropy methods, in the earthquakeMag, asteroid, and jet data sets, tend to cluster their selection more in comparison to the reconstruction-based methods. This is not to say, reconstruction-based selection methods never clump their selections. Especially, at higher budgets, (e.g. earthquakeMag with budget 20) reconstruction-based methods do make clumped selections. We merely observe that the tendency to clump selections is higher for

the joint entropy methods, providing insight into why the clustering analysis often classified them as dissimilar.

CHAPTER V

FAILURE CONDITIONS

5.1 Overview

This section mainly addresses **RQ2**: do we observe notable *outlier* behavior? Now, as mentioned in the introduction, establishing the “correct” selections for a particular combination of data set of budget is difficult without consulting a domain scientist. Therefore, this section points out *clear* cases where we observe selection methods and evaluation metrics misbehave.

5.2 Issues that data sets with high dynamic range encounter

In the context of this work, data sets with high dynamic range refers to data sets where the local minimum and maximum values of the each time step change drastically over the course of the time series. Among the data sets that considered, two data sets, `asteroid` and `earthquakeMag`, have high dynamic range. `asteroid`, for example, models an asteroid collision with a ocean that occurs at time step 4 in our data set. After this initial impact, the remaining time steps track the aftermath of the collision and how the energy disperses through the air and ocean. As a result, the range of values at each time step are drastically different. For examples the range of values (temperature in eV) at time step 4 is $[0.01870, 2.350]$, while the range of values at time step 200 is $[0.01075, 0.2057]$.

The `earthquakeMag` data set is a similar case. In this simulation of a magnitude 7.7 earthquake on the Southern San Andreas Fault, the rupture begins at time step 0, and continues until time step 54. After the fault rupture stops, the first waves reach the bounding box of the simulation by around time step 100. The waves continue to propagate the until the end of the time series. Fig. 5 shows volume rendering of the `earthquakeMag` simulation. Due to this nature of the

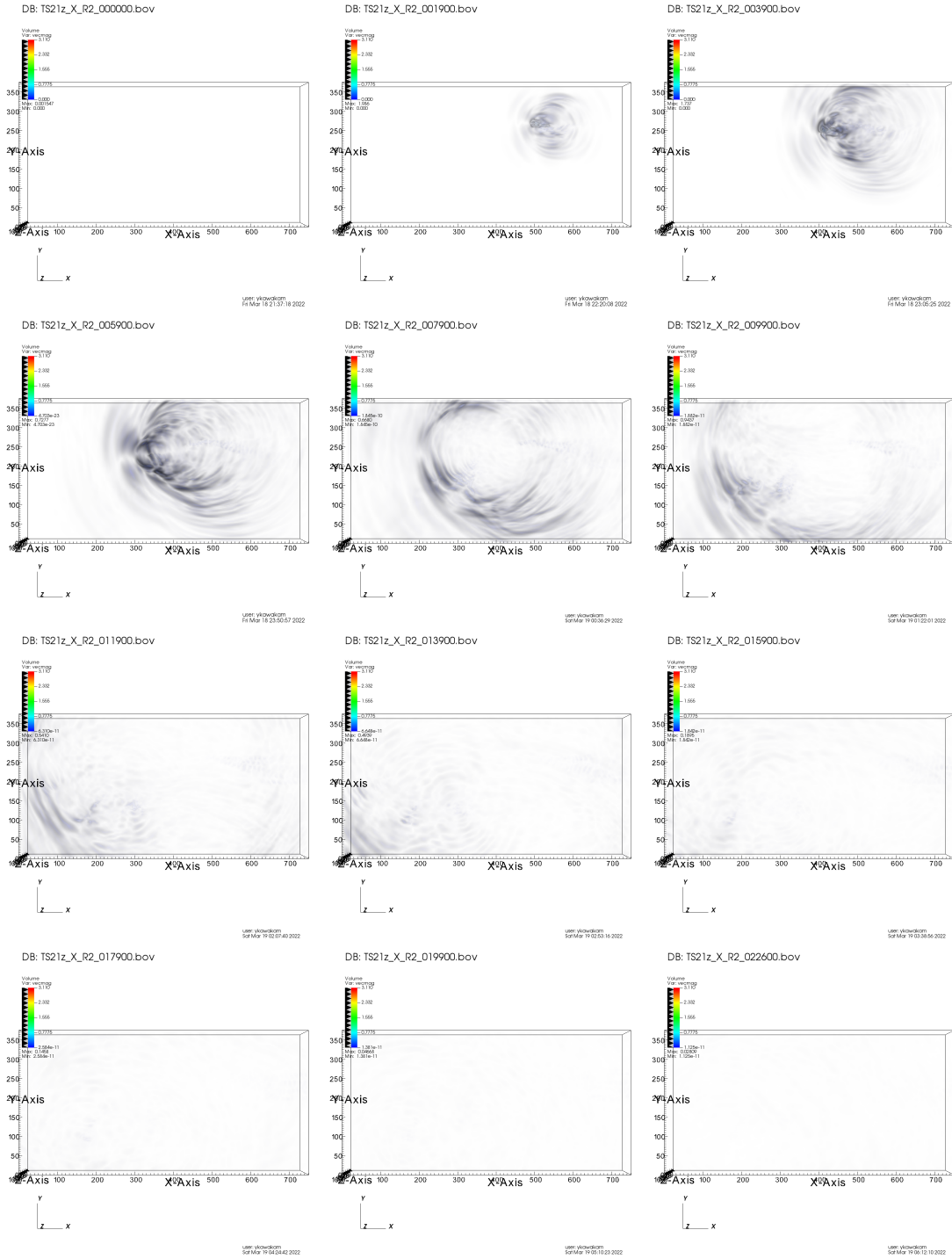


Figure 5. Volume renderings of magnitude of velocity in the earthquakeMag data set. Shown are time steps (right-to-left, top-to-bottom) 0, 19, 39, 79, 99, 119, 139, 159, 179, 199, and 226

simulation, the range of values (magnitude of velocity) at time step 39 is $[0, 3.119]$, while the range of values at time step 200 is $[0, 0.07701]$.

These drastic differences in ranges presents a challenge for any selection methods using histogram to represent data - namely the joint-entropy methods, *DTWVI* and *LIVI*. There are two options when constructing the histograms at each time step: (1) Use the local maximum and minimum values or (2) Use the global maximum and minimum values. Now, using local maximum and minimums for constructing the per time step histogram is problematic since this would effectively change resolution of each bin over the time series. Under this paradigm, entropy measures are can become practically useless since small variation of data can be distorted to represent high “information”.

As a consequence, using the global maximum and minimum values for the histograms is the only reasonable option as this ensures that the Shannon entropy numbers are meaningful in the context of the whole time series. However, using global maximum and minimum values faces one key issue. In many cases, this leads to heavy underutilization of histogram bins.

Fig. 6 shows the histogram at each time step using the global maximum and minimum values for the `earthquakeMag` data set. Note that the histograms use a log-scaled y-axis. Further, with the exception of some of the time steps that correspond to the rupture expansion and the initial wave propagation (i.e. time steps 0-54), many histogram bins are unused, or have comparatively very little count. For selection methods and evaluation metrics that rely on histograms to represent the data, this proves to be a significant challenge.

Fig. 7 plots the result from each selection method for the `earthquakeMag` data set at a budget of 10.

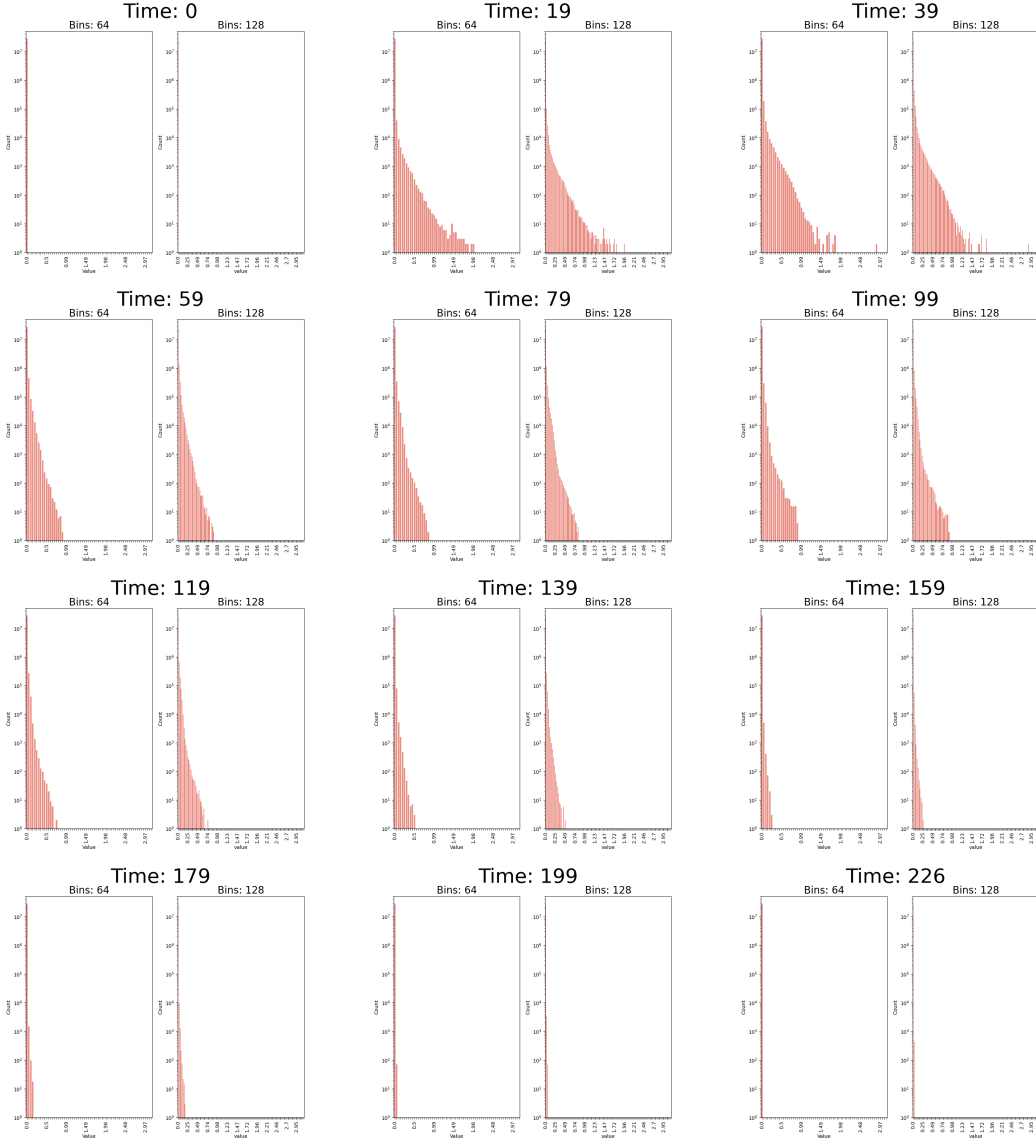


Figure 6. Histogram of time steps from the earthquakeMag data set using 64 and 128 bins. Shown are time steps (right-to-left, top-to-bottom) 0, 19, 39, 79, 99, 119, 139, 159, 179, 199, and 226.

5.2.1 Effects on joint entropy selection methods. Notice first the outlier behavior of the joint entropy approaches: JE_8 , JE_{64} , JE_{128} , and JE_{256} . While other methods maintain some temporal resolution, purely maximizing the joint entropy of the selection yields selections that are heavily clumped around where the ruptures ends and starts to spread throughout the data set.

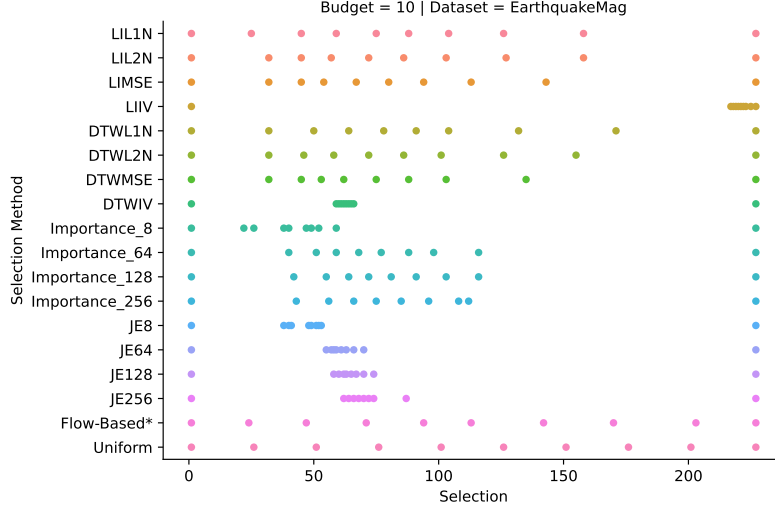


Figure 7. Result from each selection method at a budget of 10 for the earthquakeMag data set. Each dot corresponds to the time slice that each selection method makes.

(i.e. time steps 50-100) The underutilized bins in the other time steps, and the consequent empty joint histograms, imply that the joint entropy methods do not perceive these time steps as “worthy” for selection. These qualities lead to the selections by the joint entropy approaches as seen in Fig. 7. This may not be a true “failure” condition, i.e., it may be the case that this is the only region of interest to domain scientists. That said, from the perspective of understanding the simulation as a whole, this selection certainly underdelivers, and it is a notable effect of underutilized bins for data set with high dynamic range.

5.2.2 Effects on reconstruction-based selection methods

that use Variation of Information. Another notable case of the effects of underutilized bins is seen with the reconstruction-based methods that use Variation of Information as the error metric. Consider *LIVI*’s selection for the earthquakeMag data set in Fig. 7 and notice that it is clearly another outlier in relation to the rest. In fact, the selection method *LIVI*

selects $\{0, 216, 217, 218, 219, 220, 221, 222, 224, 226\}$ as its selection for the earthquakeMag data set at budget of 10. The cause of this behavior can be seen by revisiting how Variation of Information is defined.

As seen in Fig. 1, at its core, Variation of Information between two random variables X and Y is the sum of conditional entropies, $H(X|Y)$ and $H(Y|X)$. Since reconstruction-based methods try to minimize

$$\sum_i VI(data_i, data'_i)$$

where $data_i$ is the original data and $data'_i$ is the reconstructed data, it follows that this is equivalent to minimizing

$$\sum_i H(data_i | data'_i) + H(data'_i | data)$$

A key property of conditional entropy between two random variables X and Y is that it always obeys $H(X) \geq H(X|Y)$. In other words, if $H(X)$ is very small, then $H(X|Y)$ must also be a very small quantity. If X contains little information itself, then the information in X *given* the information of Y is naturally very small.

Returning to the earthquakeMag data set and its histograms in Fig. 6, notice that the first and last time step have little *information* since most data is in one bin. (i.e. $H(data_0) \approx 0$ and $H(data_{226}) \approx 0$.) This in turns means that any linearly interpolated data between time step 0 and time step 226 also must have very little *information*. (This follows since linear interpolation will never introduce values outside of what it sees at the endpoints for each mesh location.) The effect of this behavior is that, under evaluation metric EM_{LIVI} , $\{0, 226\}$ is evaluated as a better selection than $\{0, 100, 150, 226\}$. The two selections are given scores 38.041 and 55.003 respectively. (Lower is better.) The reason for this behavior is as outlined above. Since the entropy of $data'_i$ is small if linearly interpolated

between time step 0 and 226, $H(data'_i|data_i)$ also tends to be to a small amount, causing this odd effect to preferring no intermediate selections between 0 and 226. Minimization of Variation of Information can be achieved, if the reconstructed data had little to no *information* altogether. This, of course, is a clear failure condition. $\{0, 100, 150, 226\}$ as a selection contains all the information of $\{0, 226\}$; therefore, it should be evaluated at least as well as the former - not worse.

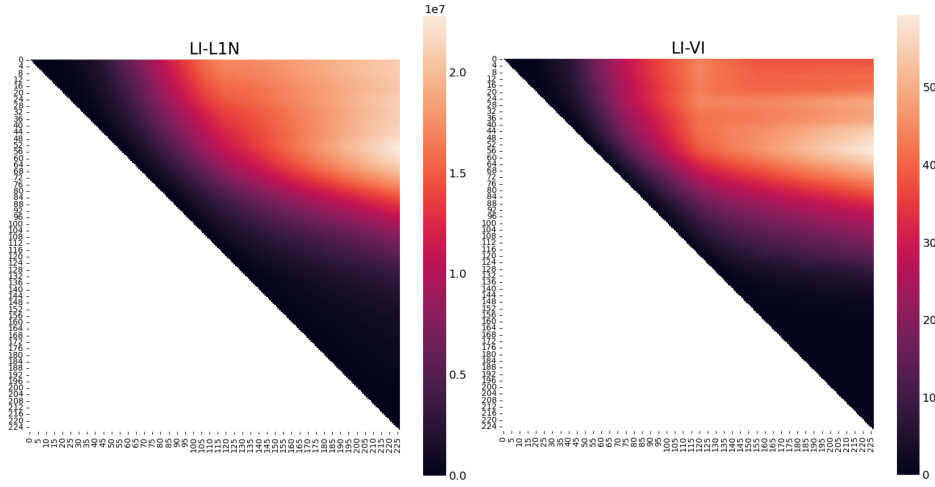


Figure 8. Pairwise mapping cost for the earthquakeMag data set: *LIL1N* on the right and *LIVI* on the left.

Fig. 8 shows the mapping cost for *LIL1N* and *LIVI* evaluation metric between any two time slices between 0 and 226 for the earthquakeMag data set. Every $[i, j]^{th}$ entry of the upper triangle heat map shows the cost between time step i and time step j under EM_{LIL1N} and EM_{LIVI} . (i.e. Linear interpolating every time step between time step i and time step j and summing the differences via *LIL1N* or *LIVI*.) The cause of outlier behavior from SM_{LIVI} can be seen in Fig. 8 as well. Notice the difference between the upper right regions of the two plots in Fig. 8. (i.e. Mapping costs between lower (0-15) time steps and higher (200-227) time steps.) While *LIL1N* assesses the comparatively high error in that region, notice

that *LIVI* does not. This provides another explanation as to why the *LIVI* selection methods makes its selection for `earthquakeMag`. Similar outlier behavior is seen with the `asteroid` data set with *LIVI* and *DTWVI*.

CHAPTER VI

METRIC COMPARISONS

6.1 Overview

This section mainly addresses **RQ3**: how do selection methods compare over many metrics? To this end, for every selection method SM in Table 3, we evaluate its selection with every evaluation metric EM in Table 3. This analysis results in the following: for every selection method SM , a budget k , and a data set, we calculate 13 metric scores, each corresponding the score that evaluation metric assigns to the selection that a particular selection method makes. However, reporting on the raw metric scores loses some key information as it relates to what the scores mean in a relative sense. For example, consider two selections S_1, S_2 , and two evaluation metrics EM_1, EM_2 . (Assume that lower is better for these two metrics.) Suppose that their evaluations are as shown in Table 5. The proper

Evaluation Metric	Evaluation of S_1	Evaluation of S_2
EM_1	10	20
EM_2	1900	2000

Table 5. Example of a set of evaluations

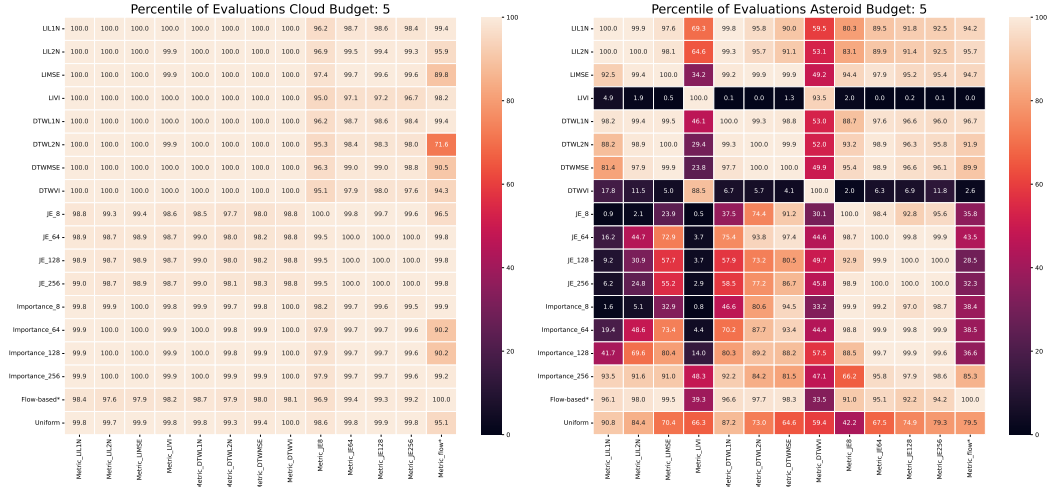
interpretation of these raw results can be difficult. For example, is the difference between 1900 and 2000 for EM_2 a meaningful difference? Or, in the eyes of EM_1 , is S_1 twice as good a selection as S_2 ? The key issue here is that raw scores from evaluation metrics lack context. Raw scores from each EM do not deliver any insight as to how each EM behaves over a range of different selections. To address this issue, for the purposes of metric comparison, this section considers instead the percentile of score among 5000 random selections rather than the raw scores itself. Specifically, when evaluating a selection S with an evaluation metric EM , we first

randomly generate 5000 selections of the same budget, evaluate each 5000 selections with EM and report the percentile of $EM(S)$ among the random 5000 selections. Note that a higher percentile is better for all metrics. For example, if a metric EM evaluates a selection S in the 75th percentile, this means that, in the eyes of EM , the selection S is better than 75% of the 5000 random selections considered.

6.2 Results

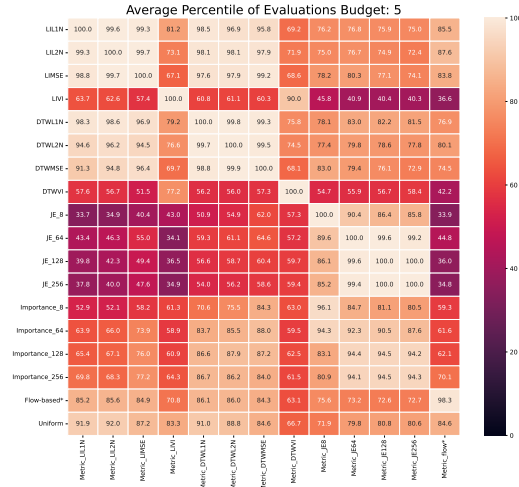
Figs. 9a and 9b plot the results from the `asteroid` and `cloud` data sets at a budget of 5 as heatmaps. Fig. 9c reports the average percentile over all 10 data sets at a budget of 5. Each row in Fig. 9 corresponds to a selection method and each column corresponds to a evaluation metric. Therefore, each cell (i, j) in the heatmaps corresponds to SM_i 's percentile among 5000 random selection under the metric EM_j .

Comparing Figs. 9a and 9b, the first obvious difference between the two is lack of poor evaluations for the `cloud` data set in Fig. 9a. Notice that every cell with the exception of $SM_{DTWL2N} - EM_{flow*}$ scores higher than a 90th percentile. In contrast, for the `asteroid` data set in Fig. 9b, the heatmap shows many more combinations with poor evaluations. This difference is primarily due to the nature of the data set. The `asteroid` data set captures a asteroid collision with a ocean, while the `cloud` data set tracks cloud over a regular day in central Europe. As a result, the `cloud` data is far more static (i.e. less changes in the data) than the `asteroid` data set. Consequently, for the `cloud` data set, all selection methods make selection very similar to an uniform selection. Since the data stays fairly similar over time, the methods fall back to placing selections mostly evenly across the time series. This is seen in Fig. 10a. The takeaway from Fig. 9a is that for



(a) cloud data set

(b) asteroid data set



(c) Average of evaluation percentiles for all data sets

Figure 9. Percentile of evaluations among 5000 random selections at a budget of 5.

static dataset like `cloud`, all methods perform similarly, and since all methods make similar selections, all evaluation methods rate them similarly as well.

On the other hand, the `asteroid` data set tells a very different story. Not only do methods make different selections (as seen in Fig. 10b), but there are notable differences in how methods fare under the range of evaluation metrics.

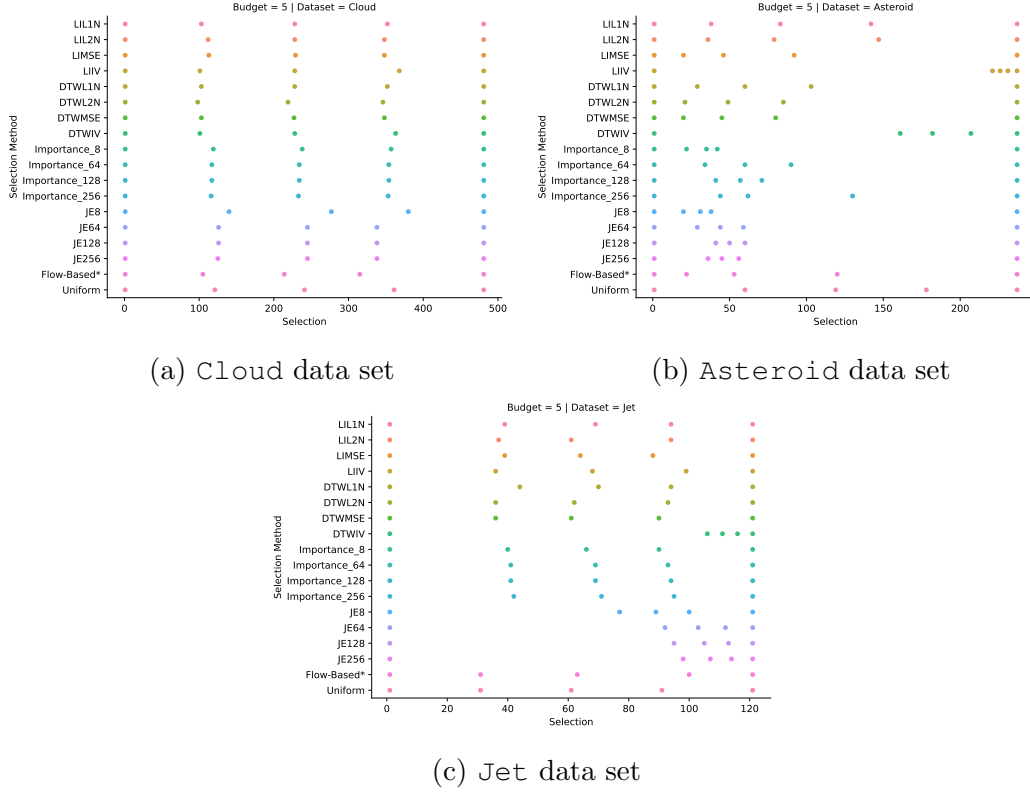


Figure 10. Results from each selection method at a budget of 5 for select data sets. Each dot corresponds to the time step that each selection method makes.

First, and perhaps most notably, notice the poor evaluation for SM_{LIVI} and SM_{DTWVI} across the board. Of course, since SM_{LIVI} and SM_{DTWVI} optimize for their respective metrics, EM_{LIVI} and EM_{DTWVI} , they perform at the 100th percentile under them. However, under any other metric, they perform poorly, around the 5-10th percentile range. The reason for this disparity becomes obvious in Fig. 10b: SM_{LIVI} and SM_{DTWVI} focus on a entirely different part of the time series than the other selection methods. The cause for this was discussed in Chapter V on failure conditions.

Second, we also notice an asymmetry between the non-VI reconstruction-based approaches (i.e. $LIL1N$, $LIL2N$, $LIMSE$, $DTWL1N$, $DTWL2N$, and $DTWMSE$) and the joint entropy approaches. While the reconstruction-based

metrics poorly evaluates the selections that the joint entropy methods make, the same is not true for the converse. In fact, the joint entropy metrics favorably evaluates selections that reconstruction-based methods make as well as the joint entropy selection methods. Notice that for the `asteroid` data set in Fig. 9b, with the exception of *LIVI* and *DTWVI*, all combinations of reconstruction-based methods and joint entropy metrics yield a higher than 80th percentile. But, several combinations of joint entropy methods and joint entropy metrics yield poor percentiles, some in the single digits. This trend holds across data sets as well, as seen in Fig. 9c. Notice that while the joint entropy metric evaluates all methods favorably, reconstruction-based metrics are more critical, especially the linear-interpolation based metrics.

Another observation we make is on the performance of uniform selection. For `cloud`, `asteroid` and `overall`, notice that uniform selection does fairly well across all metrics. In fact, in many instances, uniform selection is *preferred* by metrics over other methods that involve computations to generate. This result in both encouraging and discouraging. On the one hand, as a domain scientist, Fig. 9 shows that in many cases using a uniform selection will net better or equally as good results than applying a more sophisticated selection method. But, as a visualization researcher, this demonstrates that there is substantial room for improvement - beating a uniform selection should be a low bar considering its simplicity.

Finally, from Fig. 9c, we observe that the `flow*` metric evaluates similarly, compared to the non-VI reconstruction-based approaches. On average, the `flow*` metric favorably evaluates non-VI reconstruction-based approaches ($\sim 80^{th}$ percentile), poorly evaluates joint entropy approaches ($\sim 45^{th}$ percentile) and

evaluates the importance-based methods in between the aforementioned two.

These results are consistent with the clustering results from §4.3.3 that found the flow* selection method to make similar selections as non-VI reconstruction-based approaches.

Finally, examining the set of selection methods and their average performance over the set of evaluation metrics, we find the following

1. Non-VI reconstruction-based and flow-based* selection methods are favorably evaluated over our set of evaluation metrics, with the exception of EM_{LIVI} and EM_{DTWVI} .
2. The average performance of SM_{LIVI} and SM_{DTWVI} is poor due, in particular, to its issues outlined in Chapter V.
3. Joint entropy selection methods perform well under joint entropy evaluation metrics, but perform poorly under reconstruction-based metrics.

CHAPTER VII

CONCLUSION

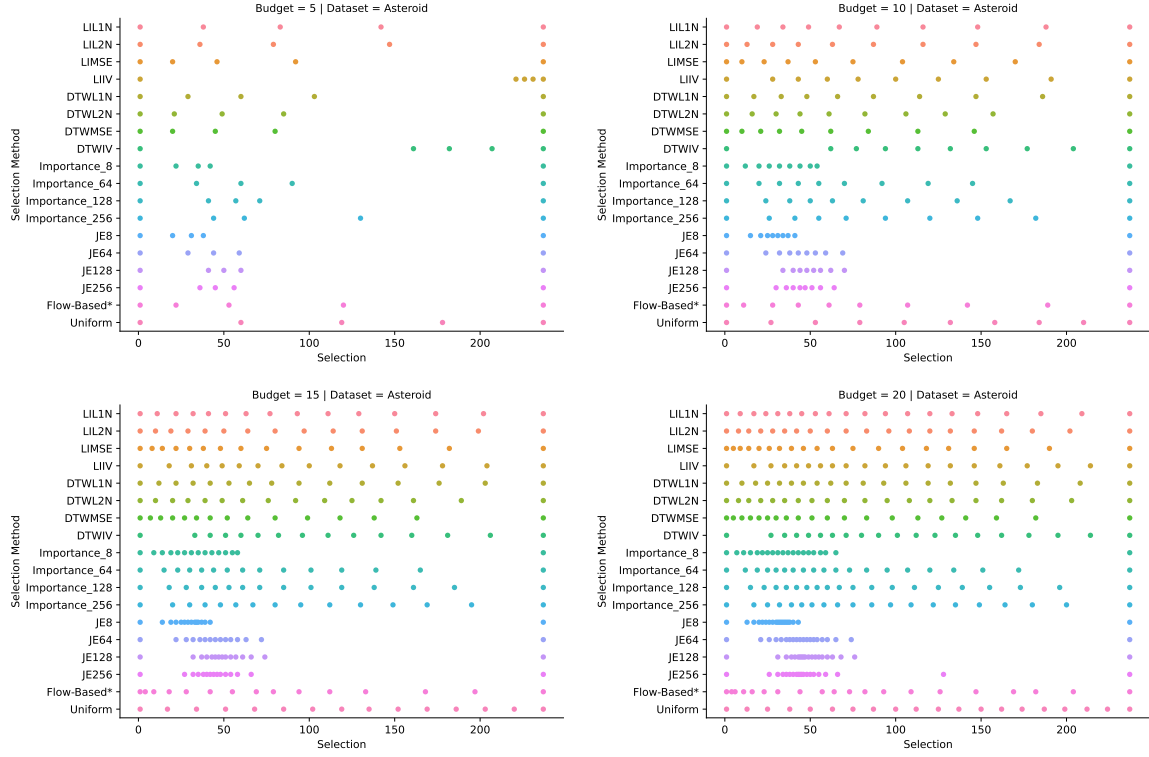
This thesis considered 10 data sets and 18 selection methods and explored the behavior of time slice selection methods over a diverse range of scientific simulations. Through a Dynamic Time Warping-alignment based selection difference metric and agglomerative clustering of selection methods, we identified similarities and dissimilarities between time slice selection algorithms in literature. Furthermore, we investigated failure conditions that some methods face. We identified that selection method that rely on histograms to represent data struggles when simulations have high dynamic range, since this leads to significantly underutilized bins in many cases. Lastly, by evaluating every selection method with every other evaluation metric, we illuminated the behavior of evaluation metrics, and identified selection methods that perform well across the set of all evaluation metrics.

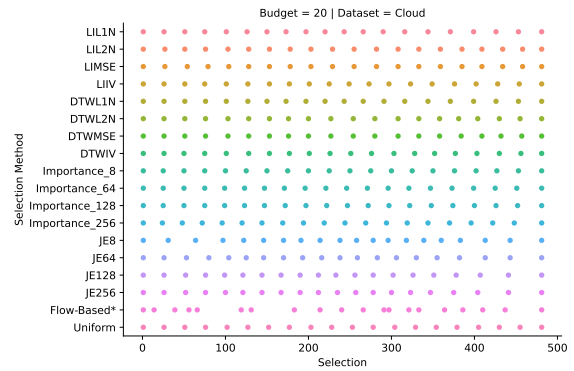
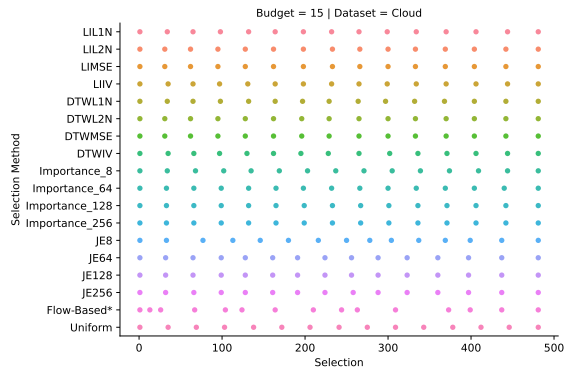
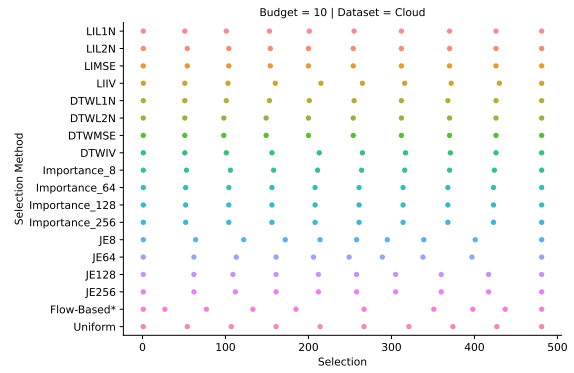
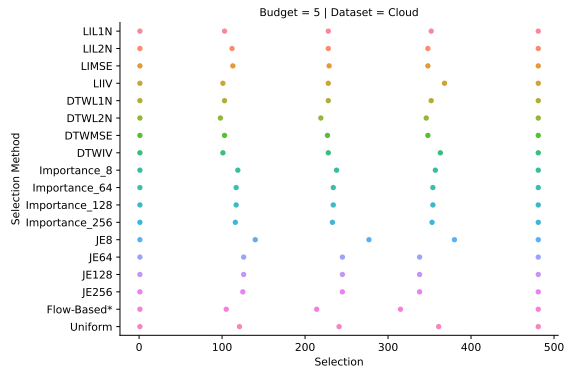
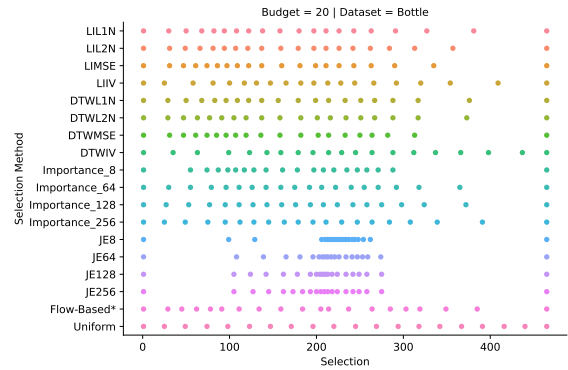
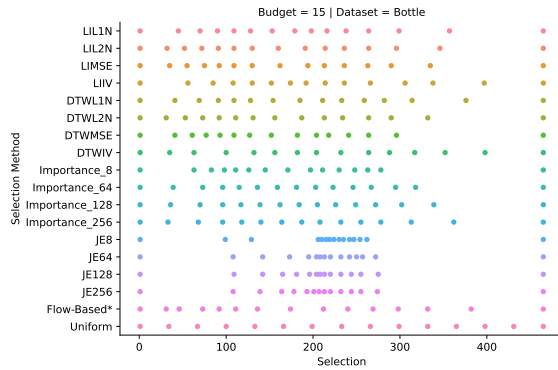
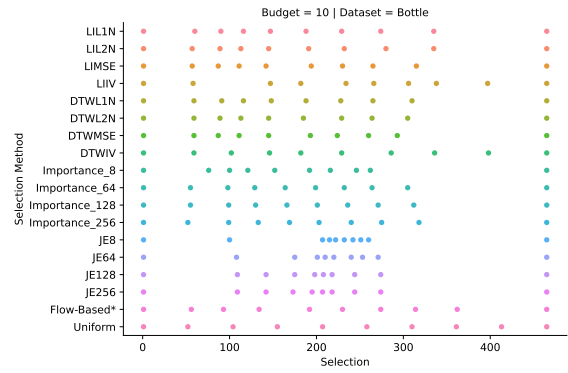
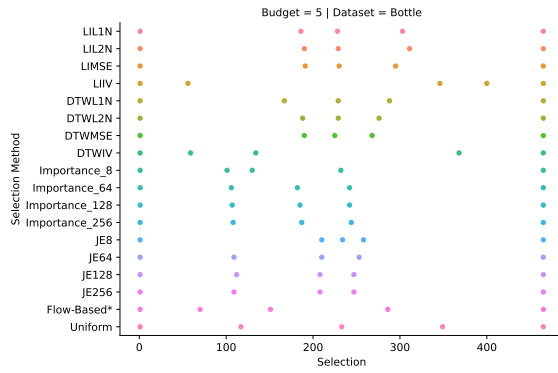
We believe this work can serve as a launching pad for further investigation into the behavior of selection methods. For example, the stability of selection methods is an important consideration. If the behavior of a selection method were significantly different for spatially subsampled data compared to the original data, this would be a cause of concern. Considering the stability of selection methods under temporal and spatial subsampling and other perturbations of the data, is an area of future work. Deep-learning based selection methods that were introduced in Chapter II may also be a worthy topic of investigation, considering the demonstrated power of deep learning in other domains of computer science. That said, this thesis has focused on domain-agnostic approaches, and deep-learning methods will need to innovate to become domain agnostic. Future work

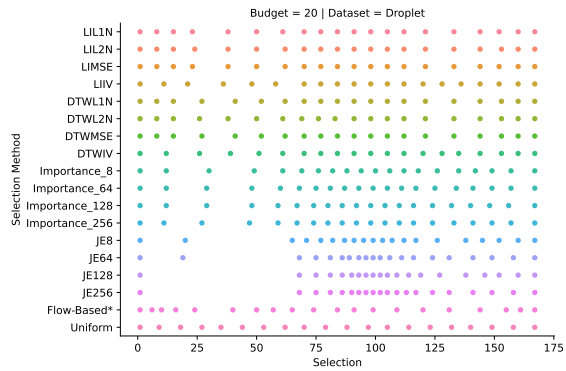
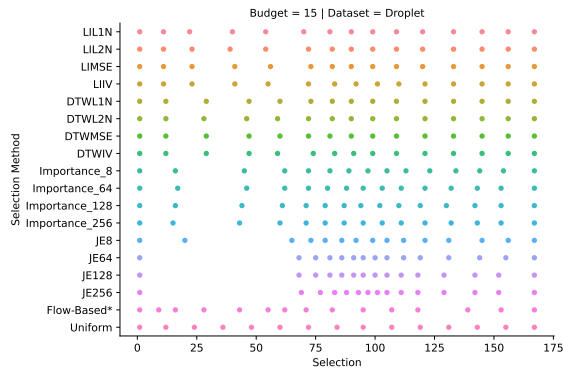
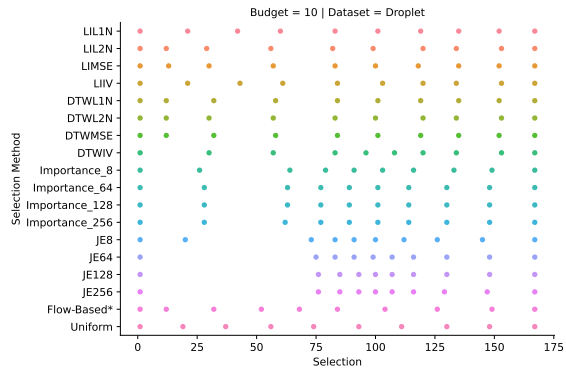
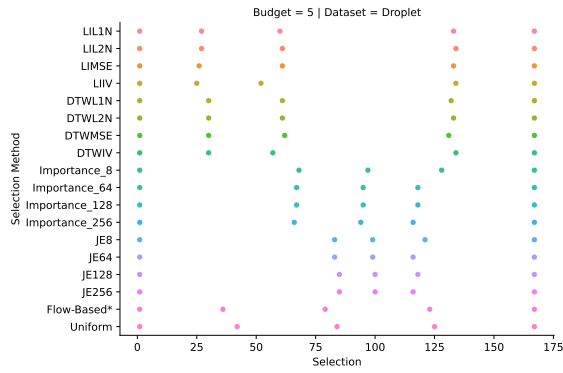
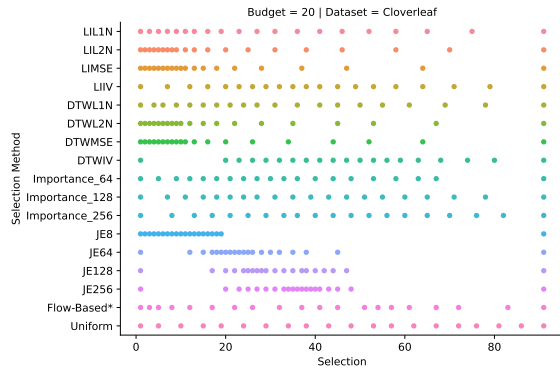
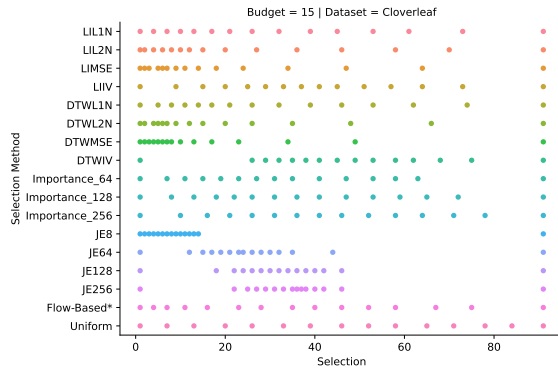
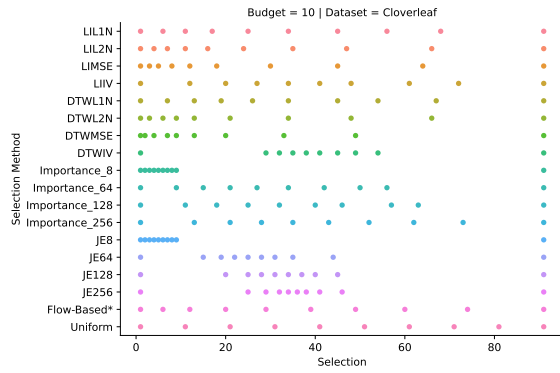
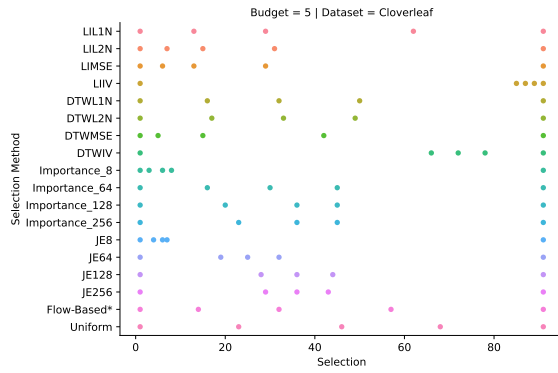
should include the aforementioned stability work, investigating deep-learning based selection methods, assessing how selection methods behave over different budgets, as well as expanding the set of considered data sets. Finally, we note that a user survey targeting domain scientists is critical to gain the full picture of the space of time slice selection methods. Recall from Chapter I that we currently consider the “correct” selection for any data set and budget to be unknown. Clearly, this is a significant obstacle in designing time slice selection methods since the real quality of selection methods can only be defined in relation to a known ground truth. In our case, the ground truth is what a domain scientist would select based on their knowledge of their domain and dynamics of a particular simulation. As such, we consider this user survey a critical future work as well.

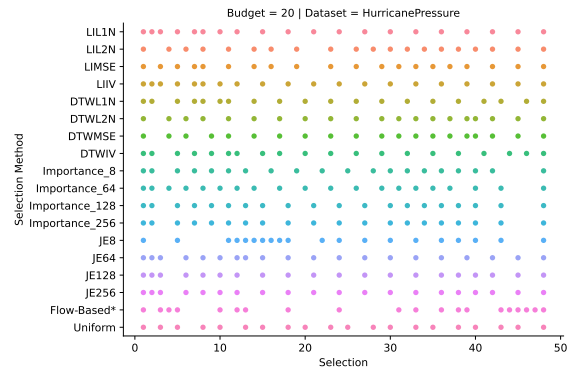
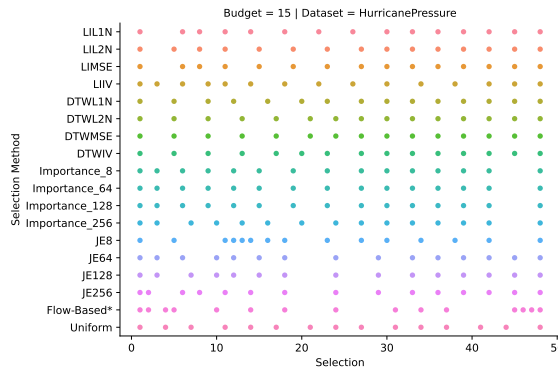
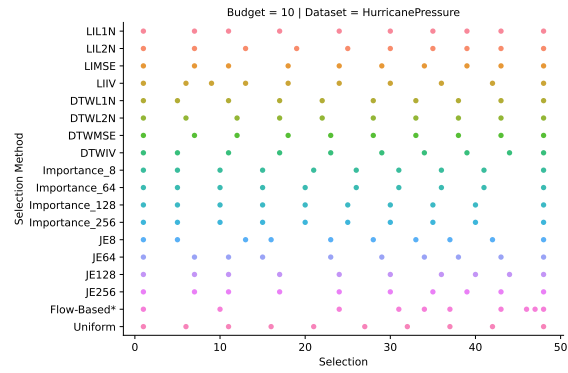
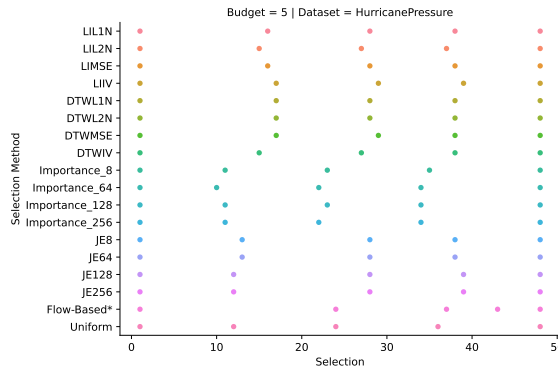
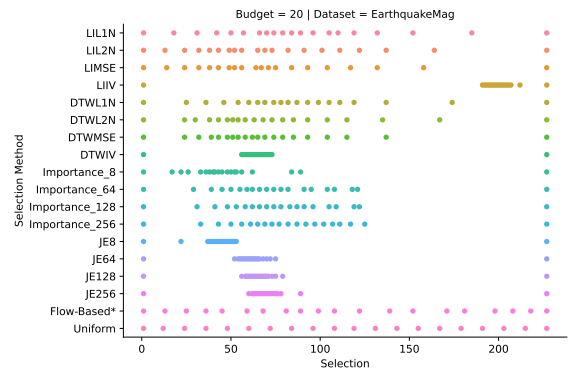
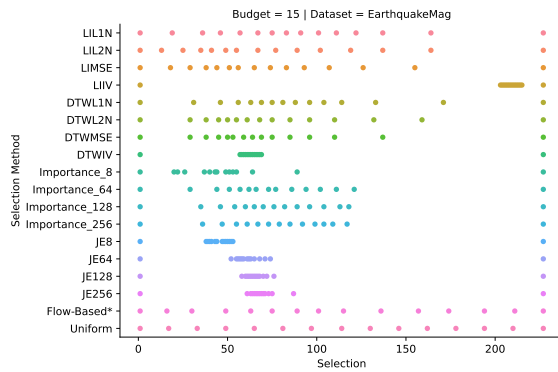
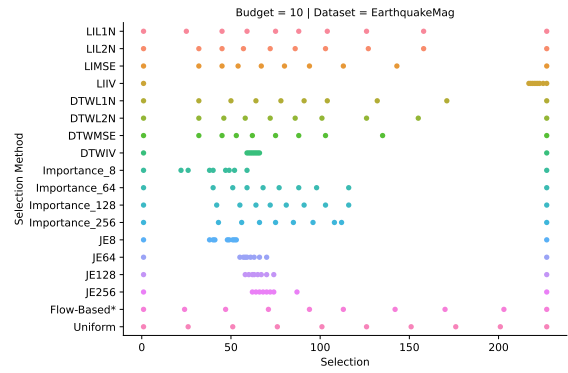
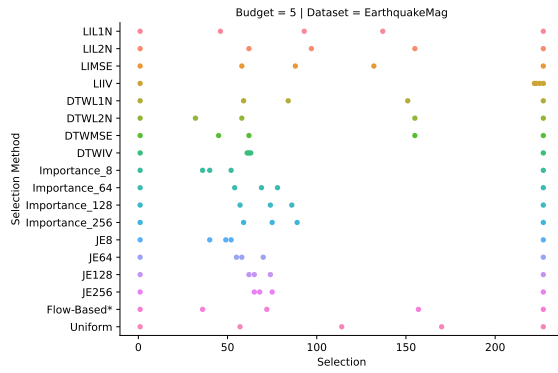
APPENDIX A

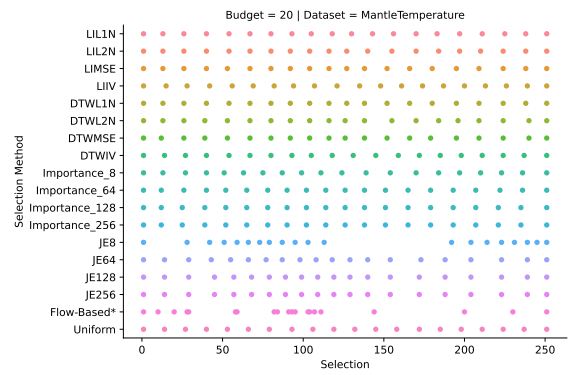
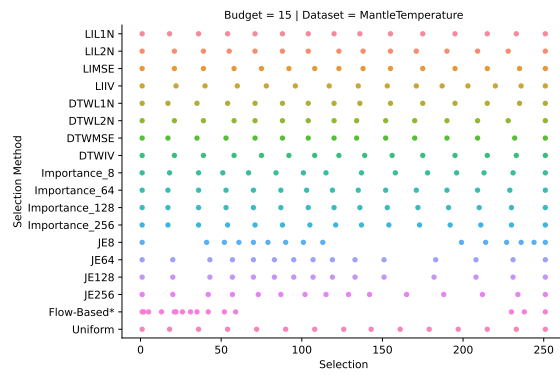
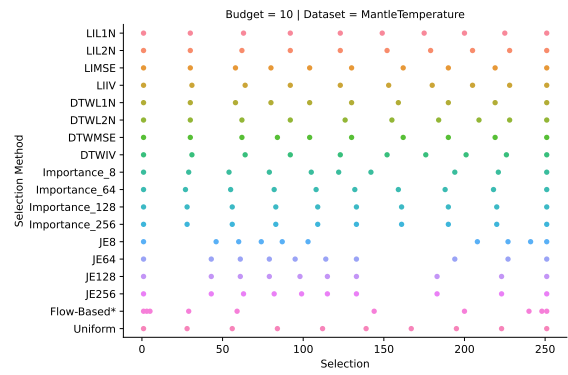
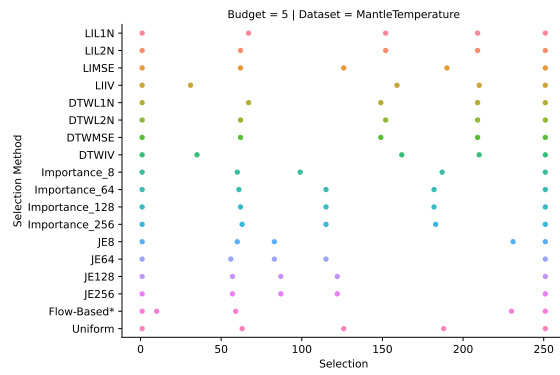
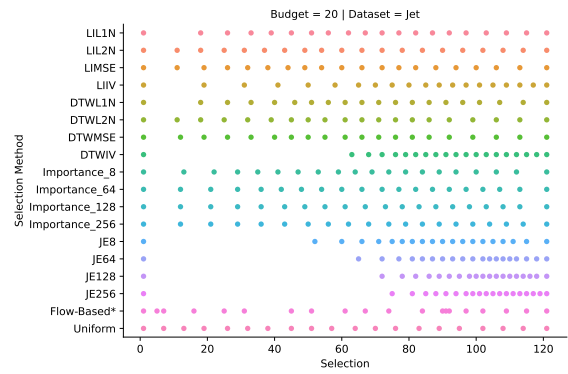
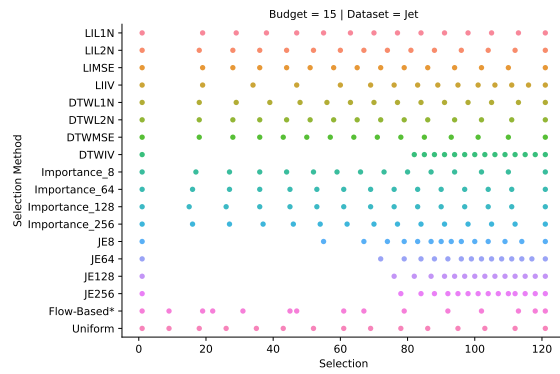
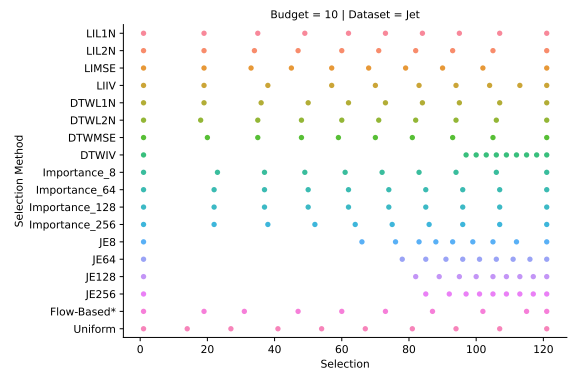
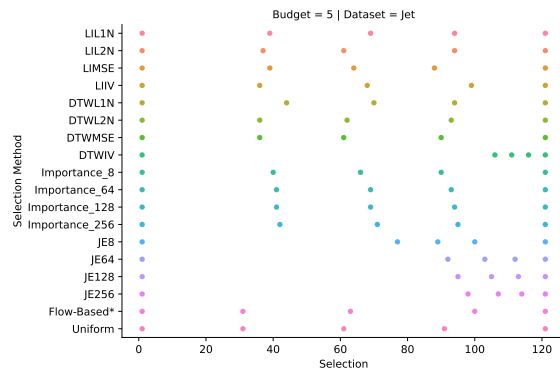
ADDITINAL RESULTS FROM EACH SELECTION METHOD

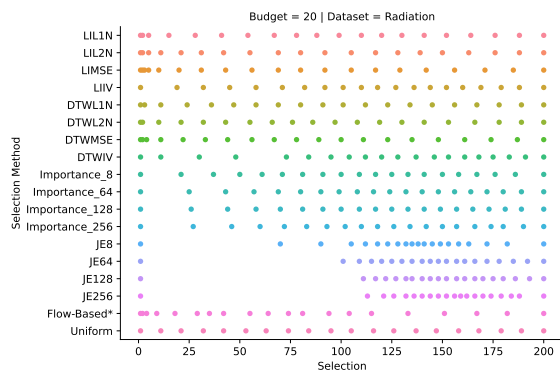
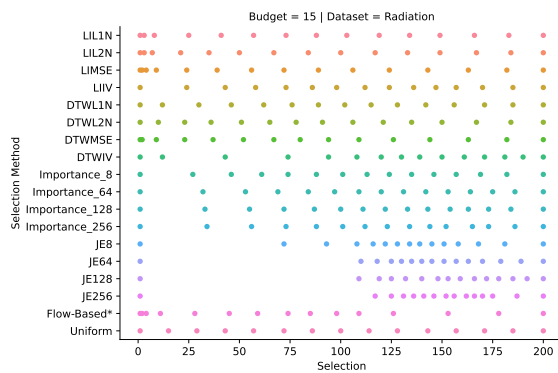
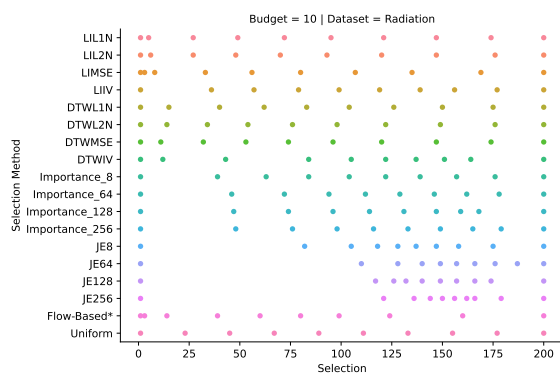
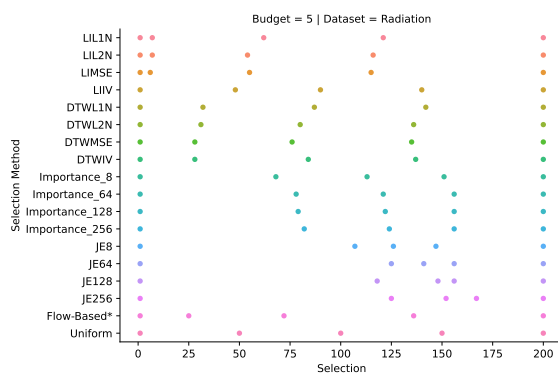








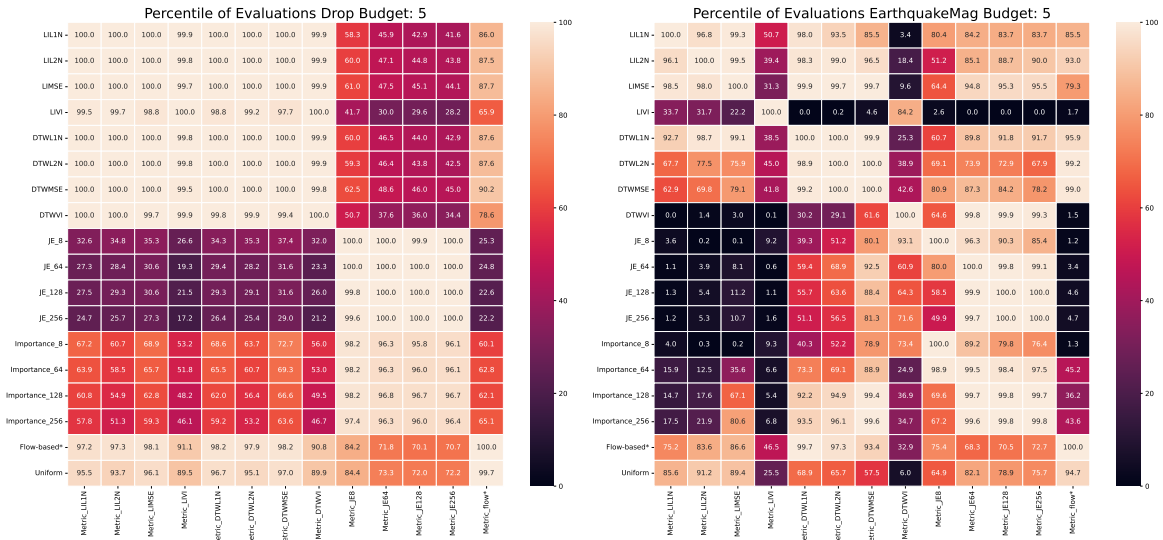


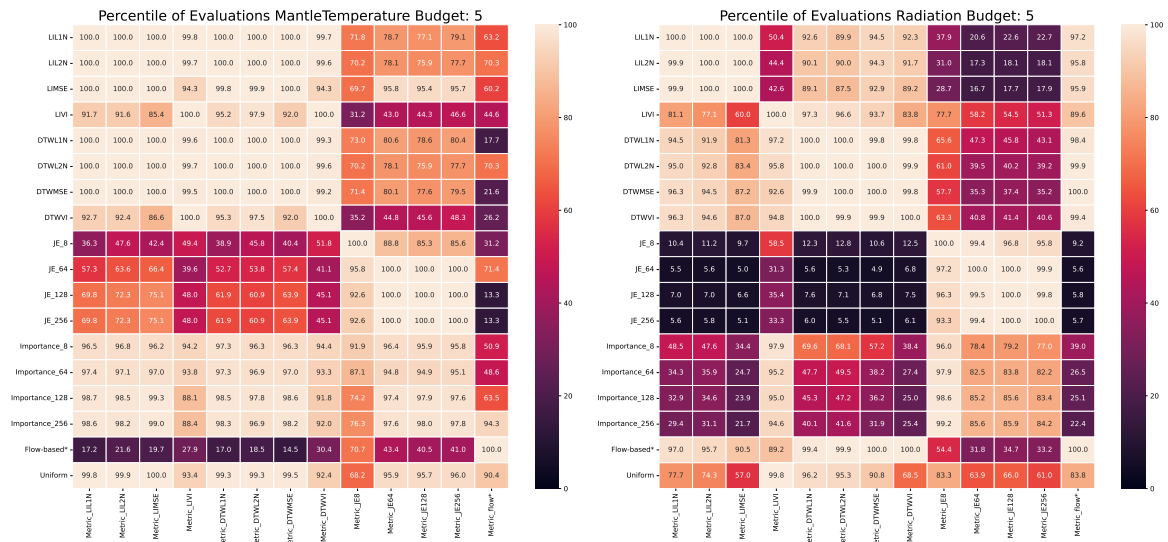


APPENDIX B

EVALUATIONS COMPARING TO RANDOM SELECTIONS







REFERENCES CITED

- [1] 2006 IEEE VISUALIZATION DESIGN CONTEST. Terashake 2.1 earthquake simulation data set.
<http://sciviscontest.ieeevis.org/2006/index.html>, 2006.
- [2] BANESH, D., WENDELBERGER, J., PETERSEN, M. R., AHRENS, J. P., AND HAMANN, B. Change point detection for ocean eddy analysis. In *EnvirVis@EuroVis* (2018), pp. 27–33.
- [3] BENNETT, J. C., BHAGATWALA, A., CHEN, J. H., PINAR, A., SALLOUM, M., AND SESHADHRI, C. Trigger detection for adaptive scientific workflows using percentile sampling. *SIAM Journal on Scientific Computing* 38, 5 (2016), S240–S263.
- [4] C. MIESTER / INSTITUTE OF AEROSPACE THERMODYNAMICS, UNIVERSITY OF STUTTGART. Droplet collision.
- [5] CHEN, M., AND JÄENICKE, H. An information-theoretic framework for visualization. *IEEE Transactions on Visualization and Computer Graphics* 16, 6 (2010), 1206–1215.
- [6] CHILDS, H., AHERN, S. D., AHRENS, J., BAUER, A. C., BENNETT, J., BETHEL, E. W., BREMER, P.-T., BRUGGER, E., COTTAM, J., DORIER, M., ET AL. A terminology for in situ visualization and analysis systems. *The International Journal of High Performance Computing Applications* 34, 6 (2020), 676–691.
- [7] CHILDS, H., BENNETT, J., GARTH, C., AND HENTSCHEL, B. In situ visualization for computational science. *IEEE Computer Graphics and Applications* 39, 6 (2019), 76–85.
- [8] DKRZ/MPI-M. Visualization of clouds and atmospheric processes.
<https://scivis2017.dkrz.de>, 2017.
- [9] FANG, Z., MOELLER, T., HAMARNEH, G., AND CELLER, A. Visualization and exploration of time-varying medical image data sets. In *Proceedings of Graphics Interface 2007* (2007), pp. 281–288.
- [10] FREY, S., AND ERTL, T. Flow-based temporal selection for interactive volume visualization. In *Computer Graphics Forum* (2017), vol. 36, Wiley Online Library, pp. 153–165.

- [11] KAWAKAMI, Y., MARSAGLIA, N., LARSEN, M., AND CHILDS, H. Benchmarking in situ triggers via reconstruction error. In *ISAV'20 In Situ Infrastructures for Enabling Extreme-Scale Analysis and Visualization*. 2020, pp. 38–43.
- [12] KEOGH, E. J., AND PAZZANI, M. J. Scaling up dynamic time warping to massive datasets. In *European Conference on Principles of Data Mining and Knowledge Discovery* (1999), Springer, pp. 1–11.
- [13] LARSEN, M., HARRISON, C., TURTON, T. L., SANE, S., BRINK, S., AND CHILDS, H. Trigger happy: Assessing the viability of trigger-based in situ analysis. In *2021 IEEE 11th Symposium on Large Data Analysis and Visualization (LDAV)* (2021), IEEE, pp. 22–31.
- [14] LING, J., KEGELMEYER, W. P., ADITYA, K., KOLLA, H., REED, K. A., SHEAD, T. M., AND DAVIS, W. L. Using feature importance metrics to detect events of interest in scientific computing applications. In *2017 IEEE 7th Symposium on Large Data Analysis and Visualization (LDAV)* (2017), IEEE, pp. 55–63.
- [15] LIU, Y., LU, Y., WANG, Y., SUN, D., DENG, L., WAN, Y., AND WANG, F. Key time steps selection for cfd data based on deep metric learning. *Computers & Fluids* 195 (2019), 104318.
- [16] MA, K.-L. In situ visualization at extreme scale: Challenges and opportunities. *IEEE Computer Graphics and Applications* 29, 6 (2009), 14–19.
- [17] MALAKAR, P., VISHWANATH, V., MUNSON, T., KNIGHT, C., HERELD, M., LEYFFER, S., AND PAPKA, M. E. Optimal scheduling of in-situ analysis for large-scale scientific simulations. In *SC'15: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis* (2015), IEEE, pp. 1–11.
- [18] MANNING, C. D., RAGHAVAN, P., AND SCHÜTZE, H. *Hierarchical clustering*. Cambridge University Press, 2018.
- [19] MYERS, K., LAWRENCE, E., FUGATE, M., BOWEN, C. M., TICKNOR, L., WOODRING, J., WENDELBERGER, J., AND AHRENS, J. Partitioning a large simulation as it runs, 2014.
- [20] NATIONAL CENTER FOR ATMOSPHERIC RESEARCH. Hurricane isabel data. <http://vis.computer.org/vis2004contest>, 2004.
- [21] PATCHETT, J. M., AND GISLER, G. R. Deep water impact ensemble data set. In *Technical Report LA-UR-17-21595* (2017).

- [22] PFISTER, H., LORENSEN, B., BAJAJ, C., KINDLMANN, G., SCHROEDER, W., AVILA, L. S., RAGHU, K., MACHIRAJU, R., AND LEE, J. The transfer function bake-off. *IEEE Computer Graphics and Applications* 21, 3 (2001), 16–22.
- [23] PORTER, W. P., XING, Y., VON OHLEN, B. R., HAN, J., AND WANG, C. A deep learning approach to selecting representative time steps for time-varying multivariate data. In *2019 IEEE Visualization Conference (VIS)* (2019), IEEE, pp. 1–5.
- [24] PULIDO, J., PATCHETT, J., BHATTARAI, M., ALEXANDROV, B., AND AHRENS, J. Selection of optimal salient time steps by non-negative tucker tensor decomposition.
- [25] PYSKLYWEC LAB/UNIVERSITY OF TORONTO. Earth’s mantle convection. <https://scivis2021.netlify.app>, 2021.
- [26] SALLOUM, M., BENNETT, J. C., PINAR, A., BHAGATWALA, A., AND CHEN, J. H. Enabling adaptive scientific workflows via trigger detection. In *Proceedings of the first workshop on in situ infrastructures for enabling extreme-scale analysis and visualization* (2015), pp. 41–45.
- [27] TONG, X., LEE, T.-Y., AND SHEN, H.-W. Salient time steps selection from large scale time-varying data sets with dynamic time warping. In *IEEE Symposium on Large Data Analysis and Visualization (LDAV)* (2012), pp. 49–56.
- [28] UK-MAC. UK-MAC/cloverleaf ref. <https://github.com/UK-MAC/CloverLeafref>.
- [29] VELTEN, A., WU, D., JARABO, A., MASIA, B., BARSİ, C., JOSHI, C., LAWSON, E., BAWENDI, M., GUTIERREZ, D., AND RASKAR, R. Femto-photography: Capturing and visualizing the propagation of light. *ACM Trans. Graph.* 32, 4 (jul 2013).
- [30] WANG, C., AND SHEN, H.-W. Information theory in scientific visualization. *Entropy* 13, 1 (2011), 254–273.
- [31] WANG, C., YU, H., AND MA, K.-L. Importance-driven time-varying data visualization. *IEEE Transactions on Visualization and Computer Graphics* 14, 6 (2008), 1547–1554.
- [32] WANG, C., YU, H., AND MA, K.-L. Importance-driven time-varying data visualization. *IEEE Transactions on Visualization and Computer Graphics* 14, 6 (2008), 1547–1554.

- [33] WANG, W., LYU, G., SHI, Y., AND LIANG, X. Time series clustering based on dynamic time warping. In *2018 IEEE 9th international conference on software engineering and service science (ICSESS)* (2018), IEEE, pp. 487–490.
- [34] WHALEN, D., AND M.L.NORMAN. Competition data set and description. <http://sciviscontest.ieeevis.org/2008/index.html>, 2008.
- [35] YAMAOKA, Y., HAYASHI, K., SAKAMOTO, N., AND NONAKA, J. In situ adaptive timestep control and visualization based on the spatio-temporal variations of the simulation results. In *Proceedings of the Workshop on In Situ Infrastructures for Enabling Extreme-Scale Analysis and Visualization* (2019), pp. 12–16.
- [36] ZHANG, Y., GUO, H., SHANG, L., WANG, D., AND PETERKA, T. A multi-branch decoder network approach to adaptive temporal data selection and reconstruction for big scientific simulation data. *IEEE Transactions on Big Data* (2021).
- [37] ZHOU, B., AND CHIANG, Y. Key time steps selection for large-scale time-varying volume datasets using an information-theoretic storyboard. *Computer Graphics Forum* 37, 3 (June 2018), 37–49.