# Compact Flow Diagrams for State Sequences

Kevin Buchin, Maike Buchin, Joachim Gudmundsson,
Michael Horton, and Stef Sijben

**Abstract.** We introduce the concept of compactly representing a large number of state sequences, e.g., sequences of activities, as a flow diagram. We argue that the flow diagram representation gives an intuitive summary that allows the user to detect patterns among large sets of state sequences. Simplified, our aim is to generate a small flow diagram that models the flow of states of all the state sequences given as input. For a small number of state sequences we present efficient algorithms to compute a minimal flow diagram. For a large number of state sequences we show that it is unlikely that efficient algorithms exist. More specifically, the problem is $W[1]$-hard if the number of state sequences is taken as a parameter. We thus introduce several heuristics for this problem. We argue about the usefulness of the flow diagram by applying the algorithms to two problems in sports analysis. We evaluate the performance of our algorithms on a football data set and generated data.

## 1   Introduction

Sensors are tracking the activity and movement of an increasing number of objects, generating large data sets in many application domains, such as sports analysis, traffic analysis and behavioural ecology. This leads to the question of how large sets of sequences of activities can be represented compactly. We introduce the concept of representing the "flow" of activities in a compact way and argue that this is helpful to detect patterns in large sets of state sequences.

To describe the problem we start by giving a simple example. Consider three objects (people) and their sequences of states, or activities, during a day. The set of state sequences $\mathcal{T} = \{\tau_1, \tau_2, \tau_3\}$ are shown in Fig. 1(a). As input we are also given a set of criteria $\mathcal{C} = \{C_1, \ldots, C_k\}$, as listed in Fig. 1(b). Each criterion is a Boolean function on a single subsequence of states, or a set of subsequences of states. For example, in the given example the criterion $C_1 =$ "eating" is true for Person 1 at time intervals 7–8am and 7–9pm, but false for all other time intervals. Thus, a criterion partitions a sequence of states into subsequences, called *segments*. In each segment the criterion is either true or false. A *segmentation* of $\mathcal{T}$ is a partition of each sequence in $\mathcal{T}$ into true segments, which is represented by the corresponding sequence of criteria. If a criterion $C$ is true for a set of subsequences, we say they *fulfil* $C$. Possible segments of $\mathcal{T}$ according to the set $\mathcal{C}$ are shown in Fig. 1(c). The aim is to summarize segmentations of all sequences efficiently; that is, build a flow diagram $\mathcal{F}$, starting at a start state $s$ and ending at an end state $t$, with a small number of nodes such that for each sequence of states $\tau_i$, $1 \leq i \leq m$, there exists a segmentation according to $\mathcal{C}$ which appears

| | Person 1 | Person 2 | Person 3 |
|---|---|---|---|
| 7-8am | breakfast | gym | breakfast |
| 8-9am | cycle to work | drive to work | cycle to work |
| 9am-5pm | work | work | work |
| 5-7pm | study | dinner | shop |
| 7-9pm | dinner | shop | dinner |

(a)

$C_1$: Eating {breakfast,dinner}
$C_2$: Commuting {cycle/drive to work}
$C_3$: Exercising {gym,cycle to work}
$C_4$: Working or studying
$C_5$: Working for at least 4 hours
$C_6$: Shopping
$C_7$: At least 2 people eating simultaneously

(b)

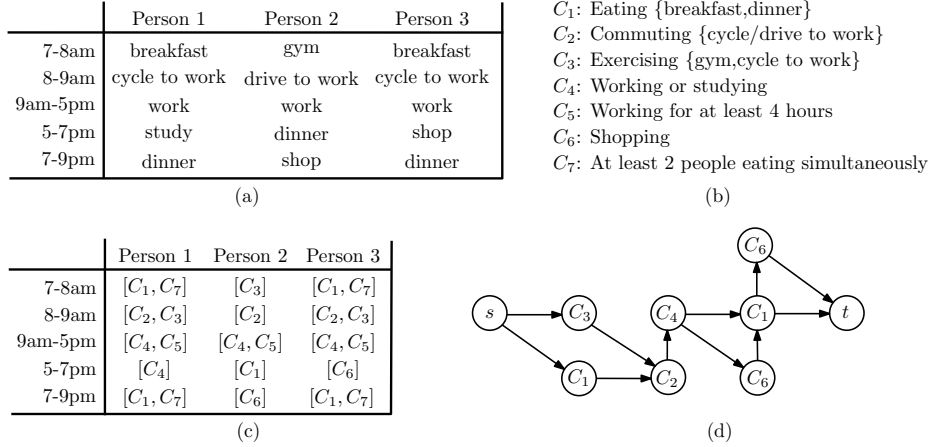| | Person 1 | Person 2 | Person 3 |
|---|---|---|---|
| 7-8am | $[C_1, C_7]$ | $[C_3]$ | $[C_1, C_7]$ |
| 8-9am | $[C_2, C_3]$ | $[C_2]$ | $[C_2, C_3]$ |
| 9am-5pm | $[C_4, C_5]$ | $[C_4, C_5]$ | $[C_4, C_5]$ |
| 5-7pm | $[C_4]$ | $[C_1]$ | $[C_6]$ |
| 7-9pm | $[C_1, C_7]$ | $[C_6]$ | $[C_1, C_7]$ |

(c)



(d)

Fig. 1: The input is (a) a set $\mathcal{T} = \{\tau_1, \ldots, \tau_m\}$ of sequences of states and (b) a set of criteria $\mathcal{C} = \{C_1, \ldots, C_k\}$. (c) The criteria partition the states into a segmentation. (d) A valid flow diagram for $\mathcal{T}$ according to $\mathcal{C}$.

as an $s$–$t$ path in $\mathcal{F}$. A possible flow diagram is shown in Fig. 1(d). This flow diagram for $\mathcal{T}$ according to $\mathcal{C}$ can be validated by going through a segmentation of each object while following a path in $\mathcal{F}$ from $s$ to $t$. For example, for Person 1 the $s$–$t$ path $s \to C_1 \to C_2 \to C_4 \to C_1 \to t$ is a valid segmentation.

Now we give a formal description of the problem. A *flow diagram* is a node-labelled DAG containing a source node $s$ and sink node $t$, and where all other nodes are labelled with a criterion. Given a set $\mathcal{T}$ of sequences of states and set of criteria $\mathcal{C}$, the goal is to construct a flow diagram with a minimum number of nodes, such that a segmentation of each sequence of states in $\mathcal{T}$ is represented, that is, included as an $s$–$t$ path, in the flow diagram. Furthermore (when criteria depend on multiple state sequences, e.g. $C_7$ in Fig. 1) we require that the segmentations represented in the flow diagram are consistent, i.e. can be jointly realized. The *Flow Diagram problem* thus requires the segmentations of each sequence of states and the minimal flow diagram of the segmentations to be computed. It can be stated as:

*Problem 1.* Flow Diagram (FD)
**Instance:** A set of sequences of states $\mathcal{T} = \{\tau_1, \ldots, \tau_m\}$, each of length at most $n$, a set of criteria $\mathcal{C} = \{C_1, \ldots, C_k\}$ and an integer $\lambda > 2$.
**Question:** Is there a flow diagram $\mathcal{F}$ with $\leq \lambda$ nodes, such that for each $\tau_i \in \mathcal{T}$, there exists a segmentation according to $\mathcal{C}$ which appears as an $s$–$t$ path in $\mathcal{F}$?

Even the small example above shows that there can be considerable space savings by representing a set of state sequences as a flow diagram. This is not a lossless representation and comes at a cost. The flow diagram represents the sequence of flow between states, however, the information about an individual sequence of states is lost. As we will argue in Section 4, paths representing many

segments in the obtained flow diagrams show interesting patterns. We will give two examples. First we consider segmenting the morphology of formations of a defensive line of football players during a match (Fig. 5). The obtained flow diagram provides an intuitive summary of these formations. The second example models attacking possessions as state sequences. The summary given by the flow diagram gives intuitive information about differences in attacking tactics.

**Properties of Criteria.** The efficiency of the algorithms will depend on properties of the criteria on which the segmentations are based. Here we consider four cases: (i) general criteria without restrictions; (ii) monotone decreasing and independent criteria; (iii) monotone decreasing and dependent criteria; and (iv) fixed criteria. To illustrate the properties we will again use the example in Fig. 1.

A criterion $C$ is *monotone decreasing* [7] for a given sequence of states $\tau$ that fulfils $C$, if all subsequences of $\tau$ also fulfil $C$. For example, if $C_4$ is fulfilled by a sequence $\tau$ then any subsequence $\tau'$ of $\tau$ will also fulfil $C_4$. This is in contrast to criterion $C_5$ which is not monotone decreasing.

A criterion $C$ is *independent* if checking whether a subsequence $\tau'$ of a sequence $\tau_i \in \mathcal{T}$ fulfils $C$ can be achieved without reference to any other sequences $\tau_j \in \mathcal{T}, i \neq j$. Conversely, $C$ is *dependent* if checking that a subsequence $\tau'$ of $\tau_i$ requires reference to other state sequences in $\mathcal{T}$. In the above example $C_4$ is an example of an independent criterion while $C_7$ is a dependent criterion since it requires that at least two objects fulfil the criterion at the same time.

**Related work.** To the best of our knowledge compactly representing sequences of states as flow diagrams has not been considered before. The only related work we are aware of comes from the area of trajectory analysis. Spatial trajectories are a special case of state sequences. A spatial trajectory describes the movement of an object through space over time, where the states are location points, which may also include additional information such as heading, speed, and temperature. For a single trajectory a common way to obtain a compact representation is *simplification* [9]. Trajectory simplification asks to determine a subset of the data that represents the trajectory well in terms of the location over time. If the focus is on characteristics other than the location, then *segmentation* [1,2,7] is used to partition a trajectory into a small number of subtrajectories, where each subtrajectory is homogeneous with respect to some characteristic. This allows a trajectory to be compactly represented as a sequence of characteristics.

For multiple trajectories other techniques apply. A large set of trajectories might contain very unrelated trajectories, hence *clustering* may be used. Clustering on complete trajectories will not represent information about interesting parts of trajectories; for this clustering on subtrajectories is needed [5,12]. A set of trajectories that forms different groups over time may be captured by a *grouping structure* [6]. These approaches also focus on location over time.

For the special case of spatial trajectories, a flow diagram can be illustrated by a simple example: trajectories of migrating geese, see [8]. The individual trajectories can be segmented into phases of activities such as directed flight,

foraging and stop overs. This results in a flow diagram containing a path for the segmentation of each trajectory. More complex criteria can be imagined that depend on a group of geese, or frequent visits to the same area, resulting in complex state sequences that are hard to analyze without computational tools.

**Organization** In Section 3 we present algorithms for the Flow Diagram problem using criteria with the properties described above. These algorithms only run in polynomial time if the number of state sequences $m$ is constant. Below we observe that this is essentially the best we can hope for by showing that the problem is $W[1]$-hard. Both theorems are proved in Section 2. Unless $W[1] = FPT$, this rules out the existence of algorithms with time complexity of $O(f(m) \cdot (nk)^c)$ for some constant $c$, where $m, n$ and $k$ are the number of state sequences, the length of the state sequences and the number of criteria, respectively. To obtain flow diagrams for larger groups of state sequences we propose two heuristics for the problem in Section 3. We experimentally evaluate the algorithms and heuristics in Section 4.

## 2 Hardness Results

In this section, the following hardness results are proven.

**Theorem 2.** *The FD problem is NP-hard. This even holds when only two criteria are used or when the length of every state sequence is* 2. *Furthermore, for any* $0 < c < 1/4$, *the FD problem cannot be approximated within factor of* $c \log m$ *in polynomial time unless* $NP \subset DTIME(m^{\mathrm{polylog}\, m})$.

Also for bounded $m$ the running times of our algorithms is rather high. Again, we can show that there are good reasons for this.

**Theorem 3.** *The FD problem parameterized in the number of state sequences is* $W[1]$-*hard even when the number of criteria is constant.*

To obtain the stated results we will perform two reductions; one from the Shortest Common Supersequence problem and one from the Set Cover problem.

### 2.1 Reduction from SCS

*Problem 4.* Shortest Common Supersequence (SCS)
**Instance:** A set of strings $R = \{r_1, r_2, \ldots, r_k\}$ over an alphabet $\Sigma$, a positive integer $\lambda$.
**Question:** Does there exist a string $s \subset \Sigma^*$ of length at most $\lambda$, that is a supersequence of each string in $R$?

The SCS problem has been extensively studied over the last 30 years (see [10] and references therein). Several hardness results are known, we will use the following two.
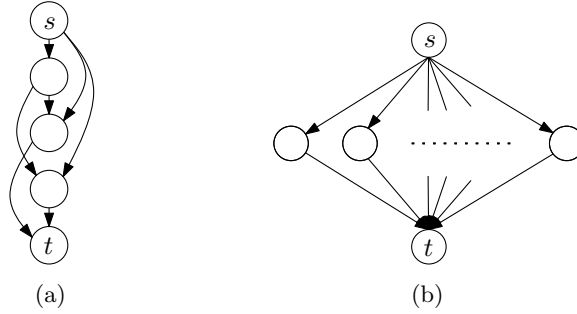
4

Fig. 2: Examples of flow diagrams produced by the reductions: (a) From SHORTEST COMMON SUPERSEQUENCE. (b) From SET COVER.

**Lemma 5 (Pietrzak [16]).** *The Shortest Common Supersequence problem parameterized in the number of Strings is $W[1]$ hard even when the alphabet has constant size.*

**Lemma 6 (Räihä and E. Ukkonen [18]).** *The Shortest Common Supersequence problem over a binary alphabet is NP-complete.*

Given an instance $I = (R = \{r_1, \ldots, r_m\}, \Sigma)$ of SCS construct an instance of FD as follows. Each character $c_l$ in the alphabet $\Sigma$ corresponds to a criterion $c_l$. Each string $r_i$ corresponds to a state sequence $T_i$, where $T_i[j] = c_{r_i[j]}$. Thus at any step $T_i$ fulfils exactly one criterion.

An algorithm for FD given an instance outputs a flow diagram $F$ of size $f$. Given $F$ one can compute a linear sequence $b$ of the vertices of $F$ using topological sort, as shown in Fig. 2a. The linear sequence $b$ has $f - 2$ vertices (omitting the start and end state of $F$) and it is a supersequence of each string in $R$. It follows that the size of $F$ is $\lambda$ if the number of characters in the SCS of $I$ has length $\lambda - 2$. Note that $F$ contains a linear sequence of vertices (after topological sort), which correspondence to a supersequence, and a set of directed edges. Consequently a solution for the FD problem can easily be transformed to a solution for the SCS problem but not vice versa.

From the above reduction, together with Lemmas 5-6, we obtain Theorem 3 and the following.

**Lemma 7.** *The FD problem is NP-hard even for two criteria.*

### 2.2 Reduction from Set Cover

*Problem 8.* Set Cover (SC)
**Instance:** A set of elements $E = \{e_1, e_2, \ldots, e_m\}$, a set of $n$ subsets of $E$, $S = \{S_1, S_2, \ldots, S_n\}$ and a positive integer $\lambda$.
**Question:** Does there exist set of $\lambda$ items in $S$ whose union equals $E$?

5

Set Cover is well known to be NP-hard, and also hard to approximate:

**Lemma 9 (Lund and Yannakakis [15]).** *For any $0 < c < 1/4$, the Set Covering problem cannot be approximated within factor of $c \log m$ in polynomial time unless $NP \subset DTIME(m^{polylog m})$.*

Given an instance $I = (E = \{e_1, e_2, \ldots, e_m\}, S = \{S_1, S_2, \ldots, S_n\})$ of Set Cover construct an instance of FD as follows. Each item $e_i$ in $E$ corresponds to a state sequence $T_i$ of length two. Each subset $S_j$ corresponds to a criterion $C_j$. If a $S_j$ contains $e_i$ then the whole state sequence $T_i$ fulfils criterion $C_j$.

An algorithm for FD given the new instance outputs a flow diagram $F$ of size $f$. The output $F$ is depicted in Fig. 2b. Given $F$ the interior vertices of $F$ corresponds to a set of subsets in $S$ whose union is $E$. The diagram $F$ has $f$ vertices if and only there is $f - 2$ subsets in $S$ that forms a Set Cover of $E$.

We obtain Theorem 2 from the above reduction, together with Lemma 9.

## 3 Algorithms

In this section, we present algorithms that compute a smallest flow diagram representing a set of $m$ state sequences of length $n$ for a set of $k$ criteria. First, we present an algorithm for the general case, followed by more efficient algorithms for the case of monotone increasing and independent criteria, the case of monotone increasing and dependent criteria, and then two heuristic algorithms.

### 3.1 General criteria

Next, we present a dynamic programming algorithm for finding a smallest flow diagram. Recall that a node $v$ in the flow diagram represents a criterion $C_j$ that is fulfilled by a contiguous segment in some of the state sequences. Let $\tau[i, j]$, $i \le j$, denote the subsequence of $\tau$ starting at the $i$th state of $\tau$ and ending at the $j$th state, where $\tau[i, i]$ is the empty sequence. Construct an $(n + 1)^m$ grid of vertices, where a vertex with coordinates $(x_1, \ldots, x_m)$, $0 \le x_1, \ldots, x_m \le n$, represents $(\tau_1[0, x_1], \ldots, \tau_m[0, x_m])$. Construct a *prefix graph* $G$ as follows:

There is an edge between two vertices $v = (x_1, \ldots, x_m)$ and $v' = (x'_1, \ldots, x'_m)$, labeled by some criterion $C_j$, if and only if, for every $i$, $1 \le i \le m$, one of the following two conditions is fulfilled: (1) $x_i = x'_i$, or (2) all remaining $\tau_i[x_i + 1, x'_i]$ jointly fulfil $C_j$. Consider the edge between $(x_1, x_2) = (1, 0)$ and $(x'_1, x'_2) = (1, 1)$ in Fig. 3(b). Here $x_1 = x'_1$ and $\tau_2[x_2 + 1, x'_2]$ fulfils $C_2$.

Finally, define $v_s$ to be the vertex in $G$ with coordinates $(0, \ldots, 0)$ and add an additional vertex $v_t$ outside the grid, which has an incoming edge from $(n, \ldots, n)$. This completes the construction of the prefix graph $G$.

Now, a path in $G$ from $v_s$ to a vertex $v$ represents a valid segmentation of some prefix of each state sequence, and defines a flow diagram that describes these segmentations in the following way: the empty path represents the flow diagram consisting only of the start node $s$. Every edge of the path adds one new node to the flow diagram, labeled by the criterion that the segments fulfil. Additionally,

6

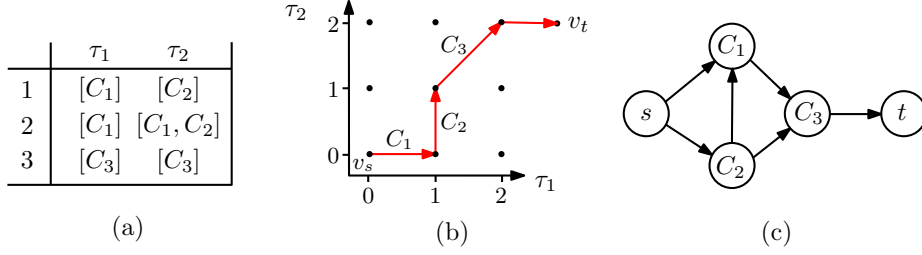| | $\tau_1$ | $\tau_2$ |
|---|---|---|
| 1 | $[C_1]$ | $[C_2]$ |
| 2 | $[C_1]$ | $[C_1, C_2]$ |
| 3 | $[C_3]$ | $[C_3]$ |

(a)     (b)     (c)

Fig. 3: (a) A segmentation of $\mathcal{T} = \{\tau_1, \tau_2\}$ according to $\mathcal{C} = \{C_1, C_2, C_3\}$. (b) The prefix graph $G$ of the segmentation, omitting all but four of the edges. (c) The resulting flow diagram generated from the highlighted path in the prefix graph.

for each node the flow diagram contains an edge from every node representing a previous segment, or from $s$ if the node is the first in a segmentation. For a path leading from $v_s$ to $v_t$, the target node $t$ is added to the flow diagram, together with its incoming edges. This ensures that the flow diagram represents valid segmentations and that each node represents at least one segment. An example of this construction is shown in Fig. 3.

Hence the length of a path (where length is the number of edges on the path) equals the number of nodes of the corresponding flow diagram, excluding $s$ and $t$. Thus, we find an optimal flow diagram by finding a shortest $v_s$–$v_t$ path in $G$.

**Lemma 10.** *A smallest flow diagram for a given set of state sequences is represented by a shortest $v_s$–$v_t$ path in $G$.*

*Proof.* We show that every $v_s$–$v_t$ path $P$ in $G$ represents a valid flow diagram $F$, with the path length equal to the flow diagram's cost, and vice versa. Thus, a shortest path represents a minimal valid flow diagram for the given state sequences.

Let $P := (v_s =: v_1, v_2, \ldots, v_\ell := v_t)$ be a $v_s$–$v_t$ path of length $\ell - 1$ in $G$. As described in the text, every $v_s$–$v_t$ path in $G$ represents a valid flow diagram, and every vertex visited by the path contributes exactly one node to the flow diagram. Thus, $P$ represents a valid flow diagram with exactly $\ell$ nodes.

For the other direction, let $F$ be a valid flow diagram of a set of state sequences $\{T_1, \ldots, T_m\}$, each of length $n$. That is, there are segmentations $\mathcal{S}_1, \ldots, \mathcal{S}_m$ of the state sequences such that every segmentation is represented in $F$ in the following way: assume the nodes of $F$ are $\{s =: f_1, f_2, \ldots, f_\ell := t\}$ according to some topological sorting. Let $\mathcal{S}_j$ consist of the segments $s_{j,1}, \ldots, s_{j,\sigma_j}$, where $\sigma_j$ is the number of segments in $\mathcal{S}_j$. Then there exists a path $(s =: f_{j,0}, f_{j,1}, \ldots, f_{j,\sigma_j}, t)$ in $F$ such that each segment $s_{j,i}$ fulfils the criterion $C(f_{j,i})$ associated with $f_{j,i}$.

Let $b_{j,i}$ be the index in $T_j$ at which $s_{j,i}$ ends, for $1 \le i \le \sigma_j$, and let $b_{j,0} := 1$. Since $\mathcal{S}_j$ is a segmentation of $T_j$, $b_{j,\sigma_j} = n$. Let $F_\lambda$ be the subdiagram of $F$ induced by $(f_1, \ldots, f_\lambda)$, for $1 \le \lambda \le \ell$. We define $x_{j,\lambda} := \max\{b_{j,i} \mid f_{j,i} \in \{f_1, \ldots, f_\lambda\}\}$. We show inductively that for each $\lambda \in \{1, 2, \ldots, \ell\}$, $G$ contains a path from $v_s$

7

to the vertex $v_\lambda := (x_{1,\lambda}, \ldots, x_{m,\lambda})$ with length $\lambda - 1$, i.e. the number of nodes in $F_\lambda$ excluding $s$.

- Base case $\lambda = 1$: Note that $v_1 = v_s$, and thus there is a path of length $\lambda - 1 = 0$ from $v_s$ to $v_1$.
- Induction step: The node $f_{\lambda+1}$ represents the segments

$$\{T_j[x_{j,\lambda}, x_{j,\lambda+1}] \mid 1 \le j \le m \land x_{j,\lambda} \ne x_{j,\lambda+1}\}.$$

Since the flow diagram is valid, these segments fulfil the criterion $C(f_{\lambda+1})$, and thus $G$ contains an edge from $v_\lambda$ to $v_{\lambda+1}$. Since a path from $v_s$ to $v_\lambda$ of length $\lambda$ exists by the induction hypothesis, there is a path from $v_s$ to $v_{\lambda+1}$ of length $\lambda + 1$.

For every state sequence $T_j$, there exists an index $\varphi_j \in \{1, \ldots, \ell - 1\}$ such that $x_{j,\lambda} = n$ for all $\lambda \ge \varphi_j$. Thus, $v_{\ell-1} = (n, n, \ldots, n)$ and $G$ contains an edge from $v_{\ell-1}$ to $v_\ell = v_t$. So, there is a path from $v_s$ to $v_t$ of length $\ell - 1$.  □

Recall that $G$ has $(n+1)^m$ vertices. Each vertex has $O(k(n+1)^m)$ outgoing edges, thus, $G$ has $O(k(n+1)^{2m})$ edges in total. To decide if an edge is present in $G$, check if the nonempty segments the edge represents fulfil the criterion. Thus, we need to perform $O(k(n+1)^{2m})$ of these checks. There are $m$ segments of length at most $n$, and we assume the cost for checking this is $T(m,n)$. Thus, the cost of constructing $G$ is $O(k(n+1)^{2m} \cdot T(m,n))$, and finding the shortest path requires $O(k(n+1)^{2m})$ time.

**Theorem 11.** *The algorithm described above computes a smallest flow diagram for a set of $m$ state sequences, each of length at most $n$, and $k$ criteria in $O((n+1)^{2m} k \cdot T(m,n))$ time, where $T(m,n)$ is the time required to check if a set of $m$ subsequences of length at most $n$ fulfils a criterion.*

### 3.2  Monotone decreasing and independent criteria

If all criteria are decreasing monotone and independent, we can use ideas similar to those presented in [7] to avoid constructing the full graph. From a given vertex with coordinates $(x_1, \ldots, x_m)$, we can greedily move as far as possible along the sequences, since the monotonicity guarantees that this never leads to a solution that is worse than one that represents shorter segments. For a given criterion $C_j$, we can compute for each $\tau_i$ independently the maximum $x_i'$ such that $\tau_i[x_i + 1, x_i']$ fulfils $C_j$. This produces coordinates $(x_1', \ldots, x_m')$ for a new vertex, which is the optimal next vertex using $C_j$. By considering all criteria we obtain $k$ new vertices. However, unlike the case with a single state sequence, there is not necessarily one vertex that is better than all others (i.e. largest ending position), since there is no total order on the vertices. Instead, we consider all vertices that are not dominated by another vertex, where a vertex $p$ dominates a vertex $p'$ if each coordinate of $p$ is at least as large as the corresponding coordinate of $p'$, and at least one of $p$'s coordinates is larger.

Let $V_i$ be the set of vertices of $G$ that are reachable from $v_s$ in exactly $i$ steps, and define $M(V) := \{v \in V \mid$ no vertex $u \in V$ dominates $v\}$ to be the set of *maximal vertices* of a vertex set $V$. Then a shortest $v_s$–$v_t$ path through $G$ can be computed by iteratively computing $M(V_i)$ for increasing $i$, until a value of $i$ is found for which $v_t \in M(V_i)$. Observe that $|M(V)| = O((n+1)^{m-1})$ for any set $V$ of vertices in the graph. Also note that $V_0 = M(V_0) = v_s$.

**Lemma 12.** *For each $i \in \{1, \ldots, \ell - 1\}$, every vertex in $M(V_i)$ is reachable in one step from a vertex in $M(V_{i-1})$. Here, $\ell$ is the distance from $v_s$ to $v_t$.*

*Proof.* Assume there exists a vertex $v \in M(V_i)$ that has no edge from a vertex in $M(V_{i-1})$. Since $v \in M(V_i)$, $v$ is also contained in $V_i$, and thus its distance from $v_s$ is $i$. Thus, there must be a vertex $v'$ at distance $i-1$ from $v_s$, i.e. $v' \in V_{i-1}$, that has an edge to $v$ representing a criterion $C_j$. By assumption, $v'$ is not contained in $M(V_{i-1})$, and thus there is a vertex $v'' \in M(V_{i-1})$ that dominates $v'$. But then, by the monotonicity of $C_j$, there must be a vertex reachable from $v''$ that is identical to $v$ or dominates $v$. Both cases lead to a contradiction. $\square$

$M(V_i)$ is computed by computing the farthest reachable vertex for each $v \in M(V_{i-1})$ and criterion, thus yielding a set $D_i$ of $O((n+1)^{m-1}k)$ vertices. This set contains $M(V_i)$ by Lemma 12, so we now need to remove all vertices that are dominated by some other vertex in the set to obtain $M(V_i)$.

We find $M(V_i)$ using a copy of $G$. Each vertex may be marked as being in $D_i$ or dominated by a vertex in $D_i$. We process the vertices of $D_i$ in arbitrary order. For a vertex $v$, if it is not yet marked, we mark it as being in $D_i$. When a vertex is newly marked, we mark its $\leq m$ immediate neighbours dominated by it as being dominated. After processing all vertices, the grid is scanned for the vertices still marked as being in $D_i$. These vertices are exactly $M(V_i)$.

When computing $M(V_i)$, $O((n+1)^{m-1}k)$ vertices need to be considered, and the maximum distance from $v_s$ to $v_t$ is $m(n+1)$, so the algorithm considers $O(mk(n+1)^m)$ vertices. We improve this bound by a factor $m$ using the following:

**Lemma 13.** *The total size of all $D_i$, for $0 \leq i \leq \ell - 1$, is $O(k(n+1)^m)$.*

*Proof.* If a vertex $v$ appears in $M(V_i)$ for some $i \in \{0, \ldots, \ell - 1\}$, it generates vertices for $D_{i+1}$ that dominate $v$, and thus $v \notin M(V_{i+j})$ for any $j > 0$. So, each of the $n^m$ vertices appears in at most one $M(V_i)$ and generates $k$ candidate vertices for $D_{i+1}$ (not all unique). Hence the total size of all $D_i$ is $O(kn^m)$. $\square$

Using this result, we compute all $M(V_i)$ in $O((k+m)(n+1)^m)$ time, since $O(k(n+1)^m)$ vertices are marked directly, and each of the $(n+1)^m$ vertices is checked at most $m$ times when a direct successor is marked. One copy of the grid can be reused for each $M(V_i)$, since each vertex of $D_{i+1}$ dominates at least one vertex of $M(V_i)$ and is thus not yet marked while processing $D_j$ for any $j \leq i$.

Since the criteria are independent, the farthest reachable point for a given starting point and criterion can be precomputed for each state sequence separately. Using the monotonicity we can traverse each state sequence once per criterion and thus need to test only $O(nmk)$ times whether a subsequence fulfils a criterion.

**Theorem 14.** *The algorithm described above computes a smallest flow diagram for $m$ state sequences of length $n$ with $k$ independent and monotone decreasing criteria in $O(mnk \cdot T(1, n) + (k + m)(n + 1)^m)$ time, where $T(1, n)$ is the time required to check if a subsequence of length at most $n$ fulfils a criterion.*

### 3.3 Monotone decreasing and dependent criteria

For monotone decreasing and dependent criteria, we can use a similar approach to that described above, however, for a given start vertex $v$ and criterion $C$, there is not a single vertex $v'$ that dominates all vertices reachable from $v$ using this criterion. Instead there may be $\Theta((n + 1)^{m-1})$ maximal reachable vertices from $v$ for criterion $C$. The maximal vertices can be found by testing $O((n+1)^{m-1})$ vertices on or near the upper envelope of the reachable vertices in $O((n+1)^{m-1} \cdot T(m, n))$ time. Using a similar reasoning as in Lemma 13, we can show that the total size of all $D_i$ $(0 \le i \le \ell - 1)$ is $O(k(n+1)^{2m-1})$, which gives:

**Theorem 15.** *The algorithm from the previous section computes a smallest flow diagram for $m$ state sequences of length $n$ with $k$ monotone decreasing criteria in $O(k(n + 1)^{2m-1} \cdot T(m, n) + m(n + 1)^m)$ time, where $T(m, n)$ is the time required to check if a set of $m$ subsequences of length at most $n$ fulfils a criterion.*

### 3.4 Heuristics

The hardness results presented in the introduction indicate that it is unlikely that the performance of the algorithms will be acceptable in practical situations, except for very small inputs. As such, we investigated heuristics that may produce usable results that can be computed in reasonable time.

For monotone decreasing and independent criteria, the heuristics we consider are based on the observation that by limiting $V_i$, the vertices that are reachable from $v_s$ in $i$ steps, to a fixed size, the complexity of the algorithm can be controlled. Given that every path in a prefix graph represents a valid flow diagram, any path chosen in the prefix graph will be valid, though not necessarily optimal. In the worst case, a vertex that advances along a single state sequence a single time-step (i.e. advancing only one state) will be selected, and for each vertex, all $k$ criteria must be evaluated, so $O(kmn)$ vertices may be processed by the algorithm. We consider two strategies for selecting the vertices in $V_i$ to retain:
(1) For each vertex in $V_i$, determine the number of state sequences that are advanced in step $i$ and retain the top $q$ vertices [*sequence heuristic*].
(2) For each vertex in $V_i$, determine the number of time-steps that are advanced in all state sequences in step $i$ and retain the top $q$ vertices [*time-step heuristic*].

In our experiments we use $q = 1$ since any larger value would immediately give an exponential worst-case running time.

## 4 Experiments

The objectives of the experiments were twofold: to determine whether compact and useful flow diagrams could be produced in real application scenarios; and to

empirically investigate the performance of the algorithms on inputs of varying sizes. We implemented the algorithms described in Section 3 using the Python programming language. For the first objective, we considered the application of flow diagrams to practical problems in football analysis in order to evaluate their usefulness. For the second objective, the algorithms were run on generated datasets of varying sizes to investigate the impact of different parameterisations on the computation time required to produce the flow diagram and the complexity of the flow diagram produced.

## 4.1 Tactical Analysis in Football

Sports teams will apply tactics to improve their performance, and computational methods to detect, analyse and represent tactics have been the subject of several recent research efforts [4,11,14,19,20,21]. Two manifestations of team tactics are in persistent and repeated occurrence of spatial formations of players, and in *plays* – a coordinated sequence of actions by players. We posited that flow diagrams would be a useful tool for compactly representing both these manifestations, and we describe the approaches used in this section.

The input for the experiments is a database containing player trajectory and event data from four home matches of the Arsenal Football Club from the 2007/08 season, provided by Prozone Sports Limited [17]. For each player and match, there is a trajectory comprising a sequence of timestamped location points in the plane, sampled at 10 Hz and accurate to 10 cm. In addition, for each match, there is a log of all the match events, comprising the timestamp and location of each event.

**Defensive Formations.** The spatial formations of players in football matches are known to characterize a team's tactics [3], and a compact representation of how formations change over time would be a useful tool for analysis. We investigated whether a flow diagram could provide such a compact representation of the defensive formation of a team, specifically to show how the formation evolves during a phase of play. The trajectories of the four defensive players were re-sampled at one-second intervals and used to compute a sequence of formation states which were then segmented to model the formation.

The criteria were derived from those presented by Kim et al. [13]. The angles between pairs of adjacent players (along the defence line) were used to compute the formation criteria, see Fig 4. We extended this scheme to allow multiple criteria to be applied where the angle between pairs of players is close to $10°$. The reason for this was to facilitate compact results by allowing for smoothing of small variations in contiguous time-steps. The criteria applied to each state is a triple $(x_1, x_2, x_3)$, computed as follows. Given two player positions $p$ and $q$, let $\angle pq$ be the angle between $p$ and $q$ relative to the goal-line. Let $R(-1) = [-90°, -5°)$, $R(0) = (-15°, +15°)$, and $R(+1) = (+5°, +90°]$ be three angular ranges. The positions of the four defenders satisfy the criteria (and thus have the formation) $(x_1, x_2, x_3)$ if $\angle p_i p_{i+1} \in R(x_i)$ for all $i \in \{1, 2, 3\}$.
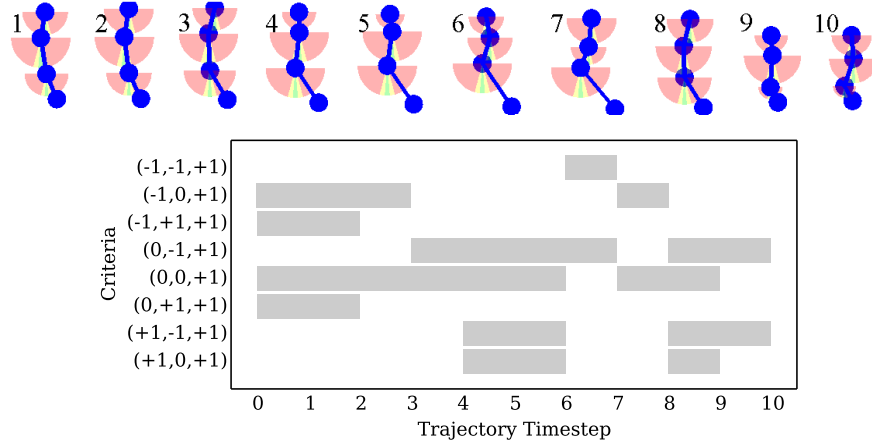
11

Fig. 4: Segmentation of a single state sequence. The formation state sequence is used to compute the segmentation representation, where segments corresponding to criteria span the state sequence *(bottom)*. The representation of this state sequence in the movement flow diagram is shaded in Fig. 5.
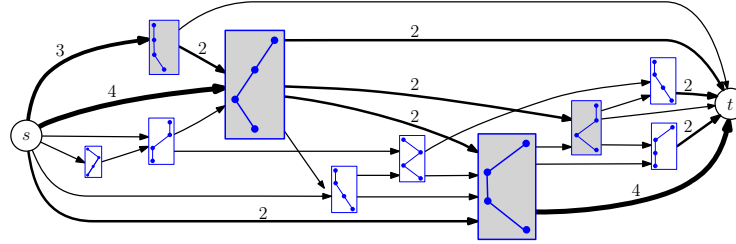


Fig. 5: Flow diagram for formation morphologies of twelve defensive possessions. The shaded nodes are the segmentation of the state sequence in Fig. 4.

The criteria in this experiment were monotone decreasing and independent, and we ran the corresponding algorithm using randomly selected sets of the state sequences as input. The size $m$ of the input was increased until the running time exceeded a threshold of 6 hours. The algorithm successfully processed up to $m = 12$ state sequences, having a total of 112 assigned segments. The resulting flow diagram, Fig. 5, has a total complexity of 12 nodes and 27 edges.

We believe that the flow diagram provides an intuitive summary of the defensive formation, and several observations are apparent. There appears to be a preference amongst the teams for the right-back to position himself in advance of the right centre-half (i.e. the third component of the triple is $+1$). Furthermore, the $(0, 0, 0)$ triple, corresponding to a "flat back four" is not present

in the diagram. This is typically considered the optimal formation for teams that utilise the offside trap, and thus may suggest that the defences here are not employing this tactic. These observations were apparent to the authors as laymen, and we would expect that a domain expert would be able to extract further useful insights from the flow diagrams.

**Attacking Plays.** During a football match, the team in possession of the ball is attempting to reach a position where they can take a shot at goal. Teams will typically use a variety of tactics to achieve such a position, e.g. teams can vary the intensity of an attack by pushing forward, moving laterally, making long passes, or retreating. We modelled attacking possessions as state sequences, segmented according to criteria representing the attacking intensity and tactics employed, and computed flow diagrams for the possessions. In particular, we were interested in determining whether differences in tactics employed by teams when playing at home or away [4] are apparent in the flow diagrams.

We focus on *ball events*, where a player touches the ball, e.g. passes, touches, dribbles, headers, and shots at goal. The event sequence for each match was divided into sub-sequences where a single team was in possession, and then filtered to include only those that end with a shot at goal.

We defined criteria that characterised the movement of the ball - relative to the goal the team is attacking - between event states in the possession sequence. The applied criteria are defined as follows. Let $x_i$, $y_i$, $t_i$ be the $x$-coordinate in metres, $y$-coordinate in metres and time-stamp in seconds, respectively, for event $i$. The velocity of the ball in the $x$-direction at event $i$, which is in the direction of the goal, is thus $x_v = (x_{i+1} - x_i)/(t_{i+i} - t_i)$ in $m/s$. The velocity $y_v$ of the ball in the $y$ direction is computed in a similar fashion, and together are used to specify the following criteria.

- *Backward movement (BM):* $x_v < 1$, a sub-sequence of passes or touches that move in a defensive direction.
- *Lateral movement (LM):* $-5 < x_v < 5$, passes or touches that move in a lateral direction.
- *Forward movement (FM):* $-1 < x_v < 12$, passes or touches that move in an attacking direction, at a velocity in the range achievable by humans, i.e. to approximately $10m/s$.
- *Fast forward movement (FFM):* $8 < x_v$, passes or touches moving in an attacking direction at a velocity generally in exceeds of maximum human velocity.
- *Long ball (LB):* a pass travelling $30m$ in the attacking direction.
- *Cross-field ball (CFB):* a pass travelling $20m$ in the cross-field direction, and that has angle in range $[80, 100]$ or $[-80, -100]$.
- *Shot resulting in goal (SG):* a successful shot resulting in a goal.
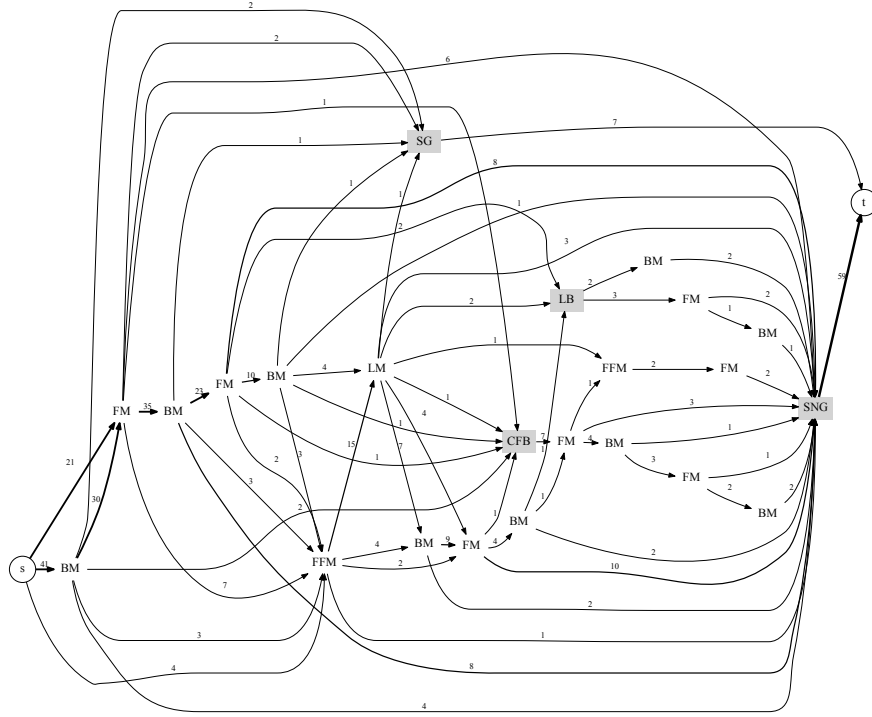- *Shot not resulting in goal (SNG):* an unsuccessful shot that does not produce a goal.

Fig. 6: Flow diagram produced for the home team. The edge weights are the number of possessions that span the edge, and the nodes with grey background are event types that are significant, as defined in Section 4.1.

For a football analyst, the first four criteria are simple movements, and are not particularly interesting. The last four events are significant: the long ball and cross-field ball change the locus of attack; and the shot criteria represent the objective of an attack.

The possession state sequences for the home and visiting teams were segmented according to the criteria and the time-step heuristic algorithm was used to compute the flow diagrams. The home-team input consisted of 66 sequences covered by a total of 866 segments, and resulted in a flow diagram with 25 nodes and 65 edges, see Fig. 6. Similarly, the visiting-team input consisted of 39 state sequences covered by 358 segments and the output flow diagram complexity was 22 nodes and 47 edges, as shown in Fig. 7.

At first glance, the differences between these flow diagrams may be difficult to appreciate, however closer inspection reveals several interesting observations. The $s$–$t$ paths in the home-team flow diagram tend to be longer than those in the visiting team's, suggesting that the home team tends to retain possession of the ball for longer, and varies the intensity of attack more often. Moreover, the
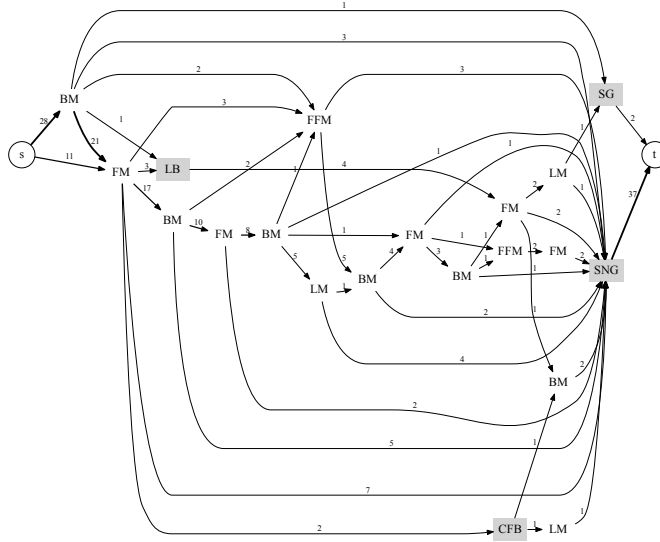
14

Fig. 7: Flow diagram produced for the visiting team. The edge weights are the number of possessions that span the edge, and the nodes with grey background are event types that are significant, as defined in Section 4.1.

nodes for cross-field passes and long-ball passes tend to occur earlier in the $s$–$t$ paths in the visiting team's flow diagram. These are both useful tactics as they alter the locus of attack, however they also carry a higher risk. This suggests that the home team is more confident in its ability to maintain possession for long attack possessions, and will only resort to such risky tactics later in a possession. Furthermore, the tactics used by the team in possession are also impacted by the defensive tactics. As Bialkowski et al [4] found, visiting teams tend to set their defence deeper, i.e. closer to the goal they are defending. When the visiting team is in possession, there is thus likely to be more space behind the home team's defensive line, and the long ball may appear to be a more appealing tactic. The observations made from these are consistent with our basic understanding of football tactics, and suggest that the flow diagrams are interpretable in this application domain.

### 4.2 Performance Testing

In the second experiment, we used a generator that outputs synthetic state sequences and segmentations, and tested the performance of the algorithms on inputs of varying sizes.

The segmentations were generated using Markov Chain Monte Carlo sampling. Nodes representing the criteria set of size $k$ were arranged in a ring and a Markov chain constructed, such that each node had a transition probability of 0.7 to

remain at the node, 0.1 to move to the adjacent node, and 0.05 to move to the node two places away. Segmentations were computed by sampling the Markov chain starting at a random node. Thus, simulated datasets of arbitrary size $m$, state sequence length $n$, criteria set size $k$ were generated.

We performed two tests on the generated segmentations. In the first, experiments were run on the four algorithms described in Section 3 with varying configurations of $m$, $n$ and $k$ to investigate the impact of input size on the algorithm's performance. The evaluation metric used was the CPU time required to generate the flow diagram for the input. In the second test, we compared the total complexity of the output flow diagram produced by the two heuristic algorithms with the baseline complexity of the flow diagram produced by the exact algorithm for monotone increasing and independent criteria.

We repeated each experiment five times with different input sequences for each trial, and the results presented are the mean values of the metrics over the trials. Limits were set such that the process was terminated if the CPU time exceeded 1 hour, or the memory required exceeded 8GB.

The results of the first test showed empirically that the exact algorithms have time and storage complexity consistent with the theoretical worst-case bounds, Fig. 8 *(top)*. The heuristic algorithms were subsequently run against larger test data sets to examine the practical limits of the input sizes, and were able to process larger input – for example, an input of $k = 128$, $m = 32$ and $n = 1024$ was tractable – although the cost is that the resulting flow diagrams were suboptimal, but correct, in terms of their total complexity.

For the second test, we investigated the complexity of the flow diagram induced by inputs of varying parameterisations when using the heuristic algorithms. The objective was to examine how close the complexity was to the optimal complexity produced using an exact algorithm. The inputs exhibited monotone decreasing and independent criteria, and thus the corresponding algorithm was used to produce the baseline. Fig. 8 *(bottom)* summarises the results for varying input parameterisations. The complexity of the flow diagrams produced by the two heuristic algorithms are broadly similar, and increase at worst linearly as the input size increases. Moreover, while the complexity is not optimal it appears to remain within a constant factor of the optimal, suggesting that the heuristic algorithms could produce usable flow diagrams for inputs where the exact algorithms are not tractable.

## 5   Concluding Remarks

We introduced flow diagrams as a compact representation of a large number of state sequences. We argued that this representation gives an intuitive summary allowing the user to detect patterns among large sets of state sequences, and gave several algorithms depending on the properties of the segmentation criteria. These algorithms only run in polynomial time if the number of state sequences $m$ is constant, which is the best we can hope for given the problem is $W[1]$-hard. As a result we considered two heuristics capable of processing large data sets in
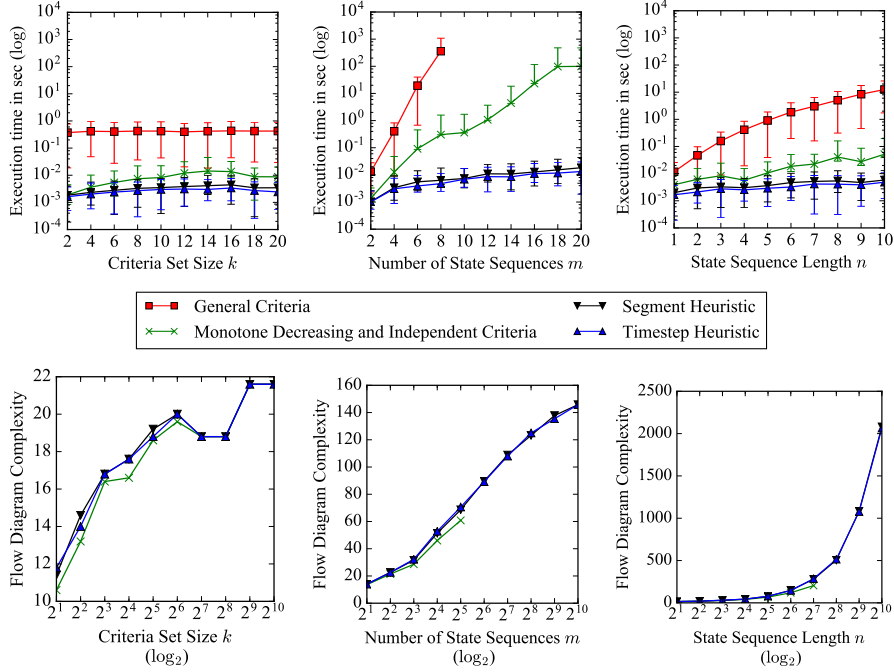
Fig. 8: Runtime statistics for generating flow diagram *(top)*, and total complexity of flow diagrams produced *(bottom)*. Default values of $m = 4$, $n = 4$ and $k = 10$ were used. The data points are the mean value and the error bars delimit the range of values over the five trials run for each input size.

reasonable time, however we were unable to give an approximation bound. We tested the algorithms experimentally to assess the utility of the flow diagram representation in a sports analysis context, and also analysed the performance of the algorithms of inputs of varying parameterisations.

## References

1. S. P. A. Alewijnse, K. Buchin, M. Buchin, A. Kölzsch, H. Kruckenberg, and M. Westenberg. A framework for trajectory segmentation by stable criteria. In *Proc. 22nd ACM SIGSPATIAL/GIS*, pages 351–360. ACM, 2014.
2. B. Aronov, A. Driemel, M. J. van Kreveld, M. Löffler, and F. Staals. Segmentation of trajectories for non-monotone criteria. In *Proc. 24th ACM-SIAM SODA*, pages 1897–1911, 2013.
3. A. Bialkowski, P. Lucey, G. P. K. Carr, Y. Yue, S. Sridharan, and I. Matthews. Identifying team style in soccer using formations learned from spatiotemporal tracking data. In *ICDM Workshops*, pages 9–14. IEEE, 2014.
4. A. Bialkowski, P. Lucey, P. Carr, Y. Yue, and I. Matthews. Win at home and draw away: automatic formation analysis highlighting the differences in home and away team behaviors. In *Proc. 8th Annual MIT Sloan Sports Analytics Conference*, 2014.

5. K. Buchin, M. Buchin, J. Gudmundsson, M. Löffler, and J. Luo. Detecting commuting patterns by clustering subtrajectories. *Int. J. Comput. Geometry Appl.*, 21(3):253–282, 2011.

6. K. Buchin, M. Buchin, M. J. van Kreveld, B. Speckmann, and F. Staals. Trajectory grouping structure. In *Proc. 13th WADS*, pages 219–230, 2013.

7. M. Buchin, A. Driemel, M. van Kreveld, and V. Sacristan. Segmenting trajectories: A framework and algorithms using spatiotemporal criteria. *Journal of Spatial Information Science*, 3:33–63, 2011.

8. M. Buchin, H. Kruckenberg, and A. Kölzsch. Segmenting trajectories based on movement states. In *Proc. 15th SDH*, pages 15–25. Springer-Verlag, 2012.

9. H. Cao, O. Wolfson, and G. Trajcevski. Spatio-temporal data reduction with deterministic error bounds. *The VLDB Journal*, 15(3):211–228, 2006.

10. C. B. Fraser and R. W. Irving. Approximation algorithms for the shortest common supersequence. *Nordic Journal of Computing*, 2(3):303–325, 1995.

11. J. Gudmundsson and T. Wolle. Football analysis using spatio-temporal tools. *Computers, Environment and Urban Systems*, 47:16–27, 2014.

12. C.-S. Han, S.-X. Jia, L. Zhang, and C.-C. Shu. Sub-trajectory clustering algorithm based on speed restriction. *Computer Engineering*, 37(7), 2011.

13. H.-C. Kim, O. Kwon, and K.-J. Li. Spatial and spatiotemporal analysis of soccer. In *Proc. 19th ACM SIGSPATIAL/GIS*, pages 385–388. ACM, 2011.

14. P. Lucey, A. Bialkowski, G. P. K. Carr, S. Morgan, I. Matthews, and Y. Sheikh. Representing and Discovering Adversarial Team Behaviors Using Player Roles. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR'13)*, pages 2706–2713, Portland, OR, jun 2013. IEEE.

15. C. Lund and M. Yannakakis. On the hardness of approximating minimization problems. In *Proc. 25th ACM STOC*, pages 286–293. ACM, 1993.

16. K. Pietrzak. On the parameterized complexity of the fixed alphabet shortest common supersequence and longest common subsequence problems. *Journal of Computer and System Sciences*, 67(4):757 – 771, 2003.

17. Prozone Sports Ltd. Prozone Sports - Our technology. `http://prozonesports.stats.com/about/technology/`, 2015.

18. K.-J. Räihä and E. Ukkonen. The shortest common supersequence problem over binary alphabet is NP-complete. *Theoretical Computer Sci.*, 16(2):187 – 198, 1981.

19. J. Van Haaren, V. Dzyuba, S. Hannosset, and J. Davis. Automatically Discovering Offensive Patterns in Soccer Match Data. In *Advances in Intelligent Data Analysis XIV - 14th International Symposium, IDA 2015*, volume 9385 of *Lecture Notes in Computer Science*, pages 286–297, Saint Etienne, oct 2015. Springer.

20. Q. Wang, H. Zhu, W. Hu, Z. Shen, and Y. Yao. Discerning Tactical Patterns for Professional Soccer Teams. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '15*, pages 2197–2206, Sydney, aug 2015. ACM Press.

21. X. Wei, L. Sha, P. Lucey, S. Morgan, and S. Sridharan. Large-Scale Analysis of Formations in Soccer. In *2013 International Conference on Digital Image Computing: Techniques and Applications (DICTA)*, pages 1–8, Hobart, nov 2013. IEEE.