# Supplemental Materials for
# CoreFlow: Extracting and Visualizing Branching Patterns from Event Sequences

## Contents

# Pseudo Code for the CoreFlow Algorithm

---

**Algorithm 1** CoreFlow algorithm

---

**input**
$S$: a collection of sequences
$n$: a tree node in the branching pattern
$m$: minimum support
$F$: an ordered list of milestone events defined by user (optional)
**procedure** CoreFlow($S, n, m, F$)
    **if** $F$ is not empty **then**
        $f \leftarrow$ first element in $F$
        add $f$ as a child of $n$
        remove $f$ from $F$
        $S_0 \leftarrow$ sequences from $S$ that do not contain $f$
        $S_1 \leftarrow$ sequences from $S$ that contain $f$
        **for** each sequence $s$ in $S_1$ **do**
            $idx \leftarrow$ index of first occurrence of $f$ in $s$
            trim $s$ from 0 to $idx$
        CoreFlow($S_0, n, m, F$)
        CoreFlow($S_1, f, m, F$)
    **if** size of $S$ ¡ m **then**
        add exit as a child of $n$
        **return**
    **else**
        $e \leftarrow$ top ranked event from $S$
        $S_0 \leftarrow$ sequences from $S$ that do not contain $e$
        $S_1 \leftarrow$ sequences from $S$ that contain $e$
        **if** size of $S_1$ ¿ m **then**
            add $e$ as a child of $n$
            **for** each sequence $s$ in $S_1$ **do**
                $idx \leftarrow$ index of first occurrence of $e$ in $s$
                trim $s$ from 0 to $idx$
            CoreFlow($S_0, n, m, F$)
            CoreFlow($S_1, e, m, F$)
        **else**
            add exit as a child of $n$
            **return**

---

# Comparing branching patterns generated using different ranking functions and datasets

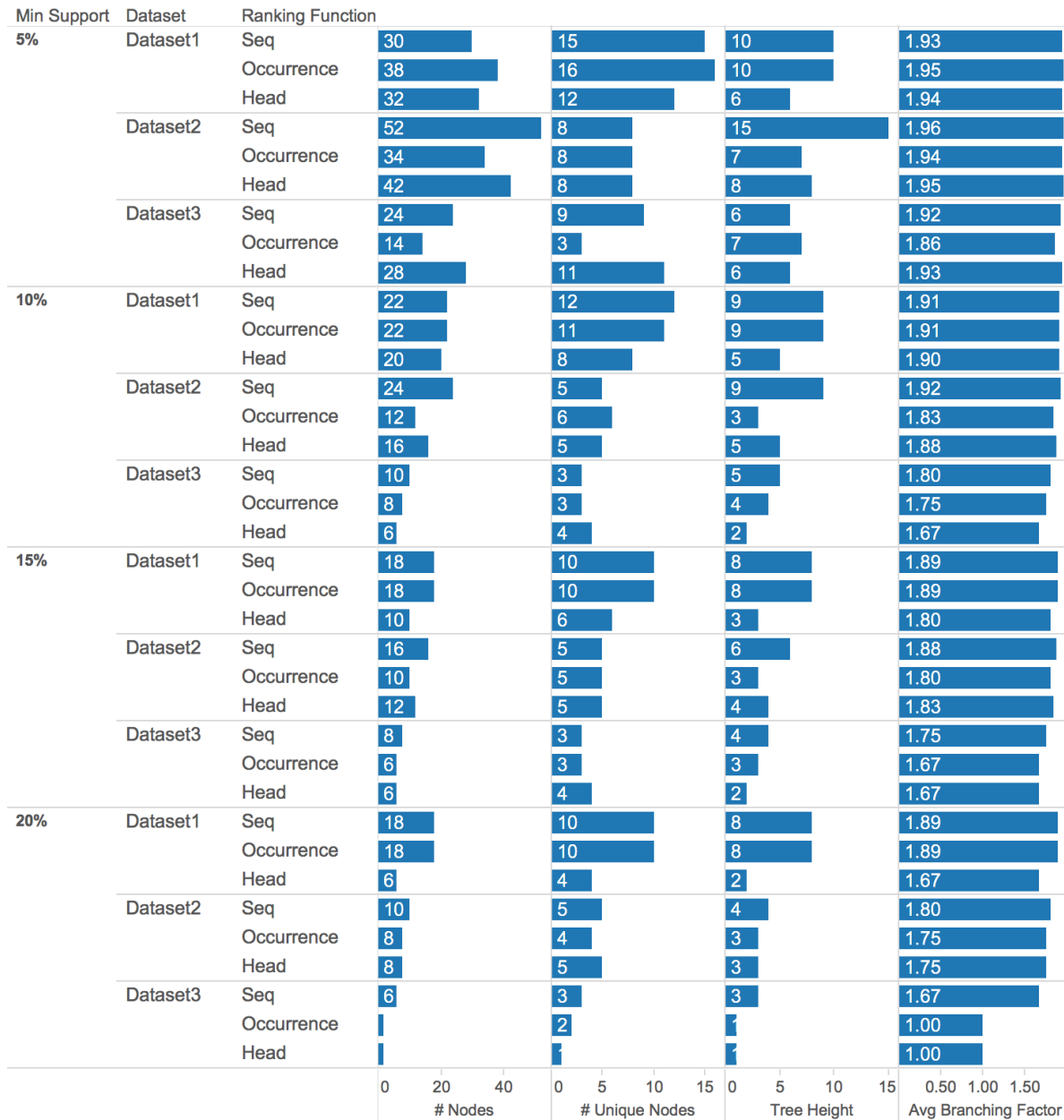| Min Support | Dataset | Ranking Function | # Nodes | # Unique Nodes | Tree Height | Avg Branching Factor |
|---|---|---|---|---|---|---|
| 5% | Dataset1 | Seq | 30 | 15 | 10 | 1.93 |
| | | Occurrence | 38 | 16 | 10 | 1.95 |
| | | Head | 32 | 12 | 6 | 1.94 |
| | Dataset2 | Seq | 52 | 8 | 15 | 1.96 |
| | | Occurrence | 34 | 8 | 7 | 1.94 |
| | | Head | 42 | 8 | 8 | 1.95 |
| | Dataset3 | Seq | 24 | 9 | 6 | 1.92 |
| | | Occurrence | 14 | 3 | 7 | 1.86 |
| | | Head | 28 | 11 | 6 | 1.93 |
| 10% | Dataset1 | Seq | 22 | 12 | 9 | 1.91 |
| | | Occurrence | 22 | 11 | 9 | 1.91 |
| | | Head | 20 | 8 | 5 | 1.90 |
| | Dataset2 | Seq | 24 | 5 | 9 | 1.92 |
| | | Occurrence | 12 | 6 | 3 | 1.83 |
| | | Head | 16 | 5 | 5 | 1.88 |
| | Dataset3 | Seq | 10 | 3 | 5 | 1.80 |
| | | Occurrence | 8 | 3 | 4 | 1.75 |
| | | Head | 6 | 4 | 2 | 1.67 |
| 15% | Dataset1 | Seq | 18 | 10 | 8 | 1.89 |
| | | Occurrence | 18 | 10 | 8 | 1.89 |
| | | Head | 10 | 6 | 3 | 1.80 |
| | Dataset2 | Seq | 16 | 5 | 6 | 1.88 |
| | | Occurrence | 10 | 5 | 3 | 1.80 |
| | | Head | 12 | 5 | 4 | 1.83 |
| | Dataset3 | Seq | 8 | 3 | 4 | 1.75 |
| | | Occurrence | 6 | 3 | 3 | 1.67 |
| | | Head | 6 | 4 | 2 | 1.67 |
| 20% | Dataset1 | Seq | 18 | 10 | 8 | 1.89 |
| | | Occurrence | 18 | 10 | 8 | 1.89 |
| | | Head | 6 | 4 | 2 | 1.67 |
| | Dataset2 | Seq | 10 | 5 | 4 | 1.80 |
| | | Occurrence | 8 | 4 | 3 | 1.75 |
| | | Head | 8 | 5 | 3 | 1.75 |
| | Dataset3 | Seq | 6 | 3 | 3 | 1.67 |
| | | Occurrence | | 2 | | 1.00 |
| | | Head | | | | 1.00 |

Figure 1: The properties of branching patterns generated using different ranking functions and datasets. We show four properties of branching patterns: number of nodes, number of unique nodes, tree height and average branching factor. The three datasets are the ones used in the evaluation study (Table 1). We use three ranking functions: 1) SEQ: the number of enclosing sequences as a metric to rank events. If a sequence has multiple occurrences of the same event, we count the sequence only once. 2) OCCURRENCE: the total number of occurrences as a metric, and 3) HEAD: the number of occurrences as the head of sequence. We also vary the minimum support from 5% to 20%. In general, HEAD tends to produce less nodes and shorter trees, while SEQ results in taller trees with more nodes.
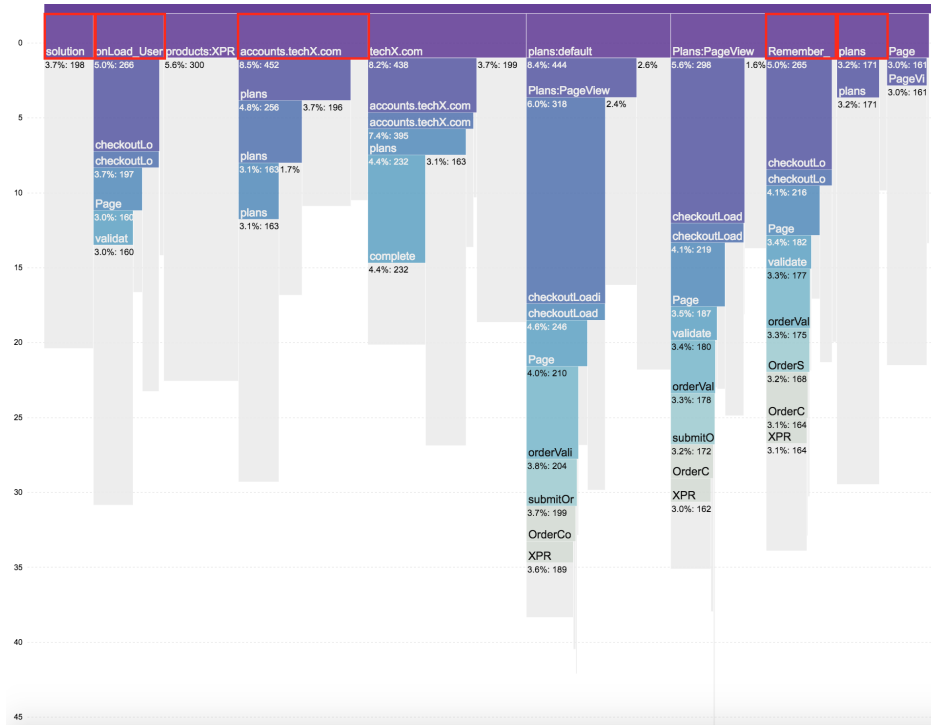
# Screenshots for Case Study 1



Figure 2: Brian grouped the sequences by the first events (entry events), and discovered that 25.4% of traffic were actually existing customers, based on the entry pages that were loaded first (these include "solution", "onLoad_User", "accounts.techX.com", "Remember_" and "plans", highlighted in red border). This insight helped him find a new user segment that he wasn't tracking already. It turned out this 'purchasers cohort' was a large driver of direct traffic for the product in general. Understanding their behavior was critical to understanding the traffic on the website.



Figure 3: Brian was able to see that of total traffic, 24.9% were switching from an existing plan (existing paid/free members) to the Single App option (highlighted in red border). This finding confirms what the company was seeing in the sales department. Many of the small and medium business customers transitioned to the main application they had been using in their full offering before.
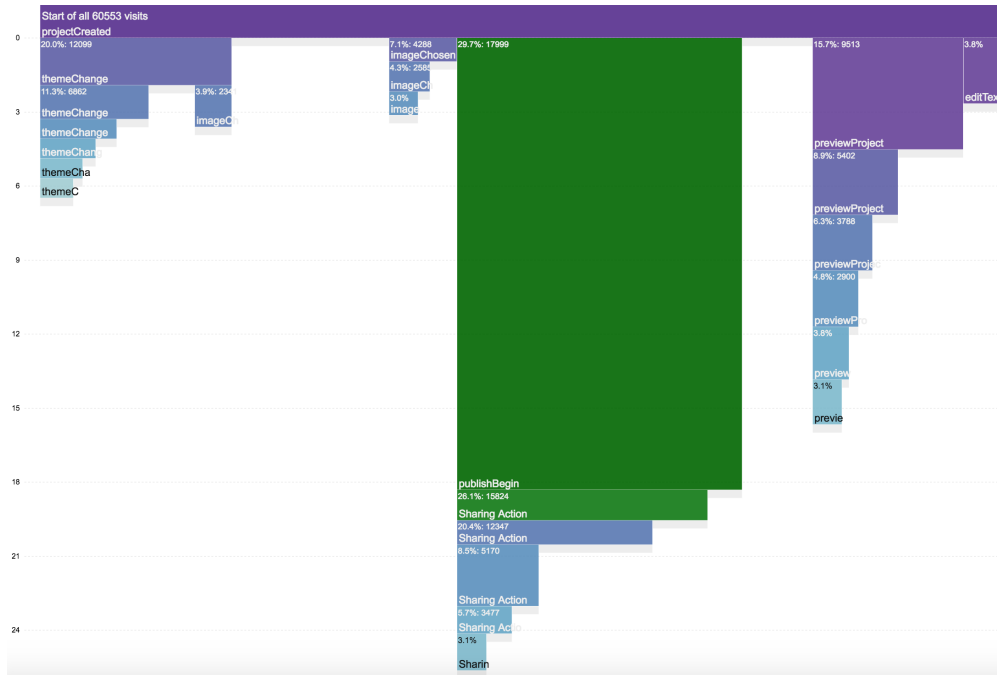
# Screenshots for Case Study 2



Figure 4: Using the funnel query, Stephanie saw that it took about 18 minutes for users to begin publishing and sharing their videos relative to the beginning of the sessions.
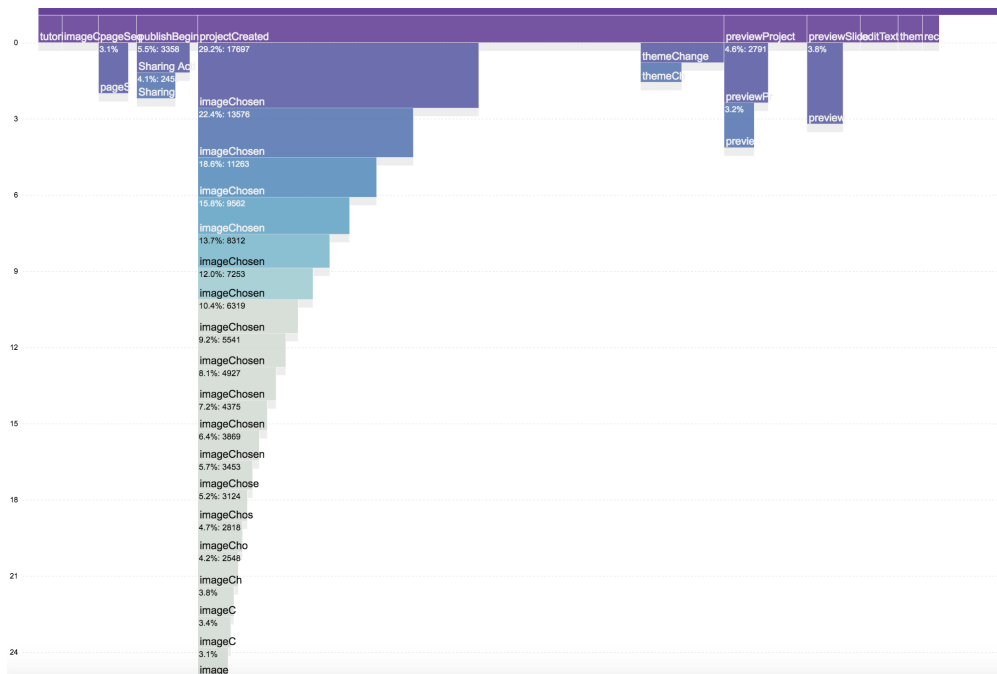


Figure 5: The most frequent user workflow was choosing images to embed in their videos.
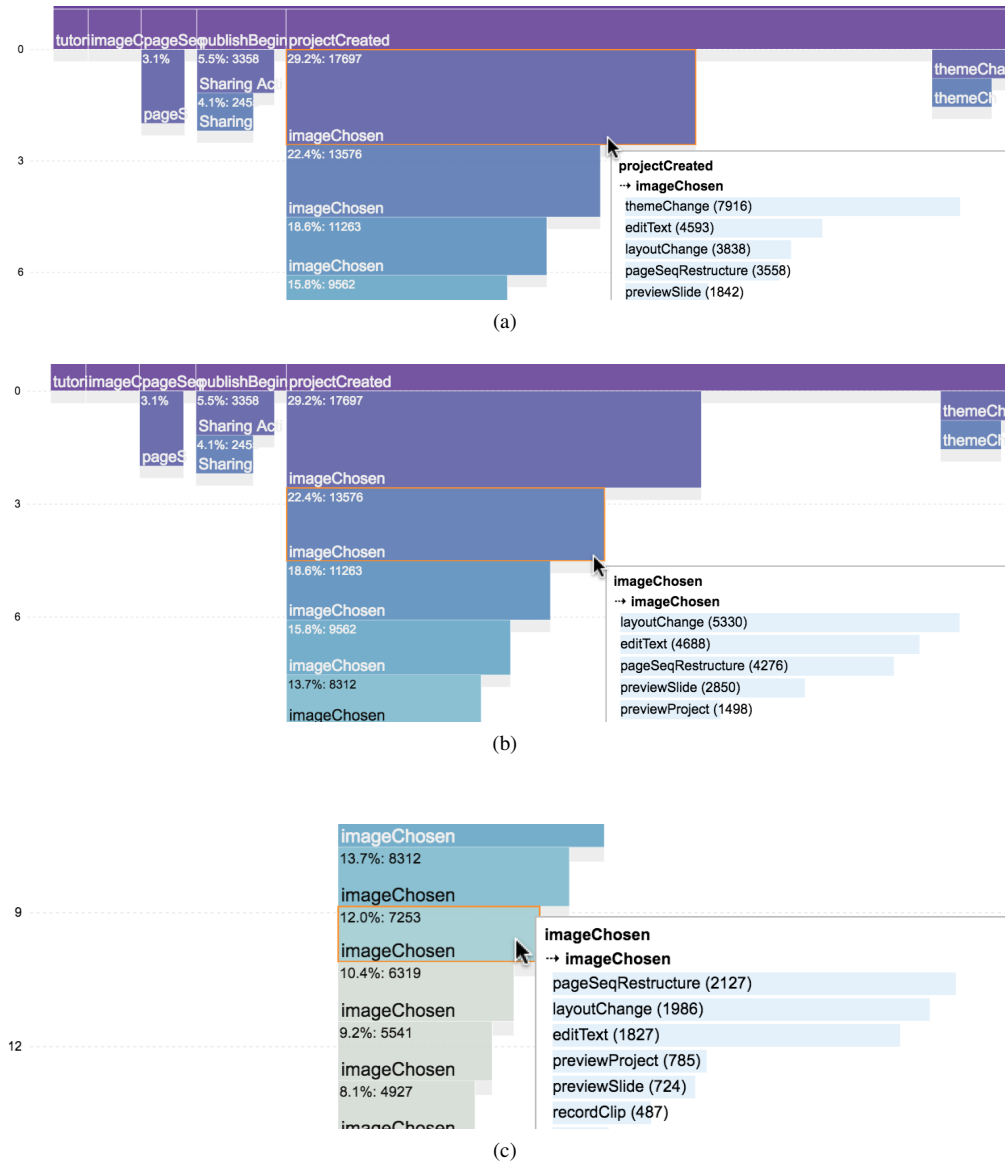
Figure 6: Users were changing video themes, changing layout and editing text amid the first few image embedding operations. As they were making progress, they started to work on page sequence restructuring more, indicating a shift of focus.