

# Towards Efficient Feedback Control in Streaming Computer Vision Pipelines

Mohamed A. Helala, Ken Q. Pu, Faisal Z. Qureshi

Faculty of Science, University of Ontario Institute of Technology  
Oshawa, ON, Canada  
{Mohamed.Helala, Ken.Pu, Faisal.Qureshi}@uoit.ca

**Abstract.** Stream processing is currently an active research direction in computer vision. This is due to the existence of many computer vision algorithms that can be expressed as a pipeline of operations, and the increasing demand for online systems that process image and video streams. Recently, a formal stream algebra has been proposed as an abstract framework that mathematically describes computer vision pipelines. The algebra defines a set of concurrent operators that can describe a pipeline of vision tasks, with image and video streams as operands. In this paper, we extend this algebra framework by developing a formal and abstract description of feedback control in computer vision pipelines. Feedback control allows vision pipelines to perform adaptive parameter selection, iterative optimization and performance tuning. We show how our extension can describe feedback control in the vision pipelines of two state-of-the-art techniques.

## 1 Introduction

Recently, there has been a rapid growth of applications capable of generating vast amounts of images and videos. Examples of such applications include online image and video sharing services (e.g., Flickr<sup>1</sup> and ImageNet<sup>2</sup>), video surveillance systems [1–4], and satellite imagery [5–7]. An emerging direction for understanding and harnessing such big visual data is *stream processing*. Stream processing represents a category of methods that process infinite sequences of data, also called data streams. In this paper, we are interested in image and video streams which we refer to as *Vision Streams*. In order to process vision streams, researchers use stream processing concepts such as pipelines to construct online computer vision algorithms [8–14]. A question that then arise is how can we formally and efficiently describe online computer vision pipelines? In database community, this question has been answered for text stream processing by developing stream algebra frameworks [15–18]. A stream algebra defines a set of abstract algebraic operators with well defined semantics that process streams as operands. These operators are used to build mathematical expressions that declaratively construct stream processing pipelines. For example, Demers *et al.* [18] proposed an algebra to express queries on event streams. Chkodrov *et al.* [15] described

<sup>1</sup> Flickr: <https://www.flickr.com/> (last accessed on 7 September 2014).

<sup>2</sup> ImageNet: <http://www.image-net.org/> (last accessed on 7 September 2014).

an implementation of a stream algebra that extends relational algebra for data streams. There are several advantages for stream algebras. For example, they can provide formal methods to resolve pipeline bottlenecks, implement dynamic re-configuration, apply incremental evaluation, define common optimization methods, etc.

Database stream algebras are designed for structured textual streams. So, they are inapplicable for vision streams with unstructured visual content. Despite such challenge, there exist some software frameworks such as OpenVL [19] and GStreamer [20], that address the efficient implementation of vision pipelines. However, these frameworks do not define a stream algebra, and lack a formal definition of vision pipelines. Recently, Helala *et al.* [21] presented a stream algebra for computer vision pipelines. This algebra revises the previous database stream algebras, and provides an abstraction for formally expressing computer vision pipelines. The algebra contains operators for both data processing and flow rate control. Two online vision algorithms have been expressed in this algebra by [21]. However, these algorithms are only for feedforward pipelines. This limits the applicability of the algebra to other online computer vision systems that use feedback control to perform tasks such as parameter tuning [22–25], and iterative optimization [14]. These tasks are widely used in online vision algorithms. For example, Sherrah [23] presented an algorithm for continuous real-time parameter tuning of a people tracking surveillance system. Supancic *et al.* [25] explored parameter tuning for long term tracking. Iterative optimization was also studied by [14] to iteratively align Flickr’s photo streams.

In this paper, we are studying feedback control loops in computer vision pipelines. Specifically, we extend the stream algebra of [21] to provide an algebraic description of feedback control; here we focus on parameter tuning, and iterative optimization. Our description of feedback control is formal and abstract, which makes it reusable by several online computer vision pipelines [22–25, 14]. The paper is organized as follows: Section 2 briefly reviews the algebra of [21], and discusses feedback control. Section 3 describes the feedback control of two state-of-the-art online computer vision algorithms. Then, we provide discussions in Section 4, and conclude the paper in Section 5.

## 2 Stream Algebra

The stream algebra in [21] contains three main parts: a common notation, a set of operators, and the formal semantics used to write pipeline expressions. This section gives a brief review of the algebra, and states the operators used in this paper. Finally, we discuss our algebraic extensions that provide an abstract and formal definition of feedback control in vision pipelines.

### 2.1 Notation

The algebra in [21] defines a data stream as an infinite sequence of data chunks with two function  $\lambda x : x \rightarrow s$ , and  $\leftarrow s$ , to write to and read from a stream

$s$ , respectively. The algebra indicates the set of all possible streams as  $\mathbf{S}$ . To indicate the type of the stream data, the notation  $\mathbf{S}\langle T \rangle$  is used, where  $T$  signify the data type. A stream operator is defined as a mapping function  $h : \mathbf{S}^m \rightarrow \mathbf{S}^n : S_{\text{in}}^1, \dots, S_{\text{in}}^m \rightarrow S_{\text{out}}^1, \dots, S_{\text{out}}^n$ , that maps  $m$  input streams to  $n$  output streams. We can define an operator by only using the following constructs as defined by [21]:

<ul style="list-style-type: none"> <li>• Shared states:               <ul style="list-style-type: none"> <li><b>state</b> <math>u</math></li> <li>-Indicates that <math>u</math> is a state for subsequent loops.</li> </ul> </li> <li>• Concurrency:               <ul style="list-style-type: none"> <li><b>loop</b> : <math>body\ of\ loop</math></li> <li>-Iterates over the body forever.</li> <li>-Each loop runs in its own thread.</li> <li>-All loops of the <i>same</i> operator share the states.</li> <li>-If there are multiple concurrent loops, <b>loop</b><sub><math>j</math></sub> indicates the <math>j</math>-th concurrent session.</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Atomicity:               <ul style="list-style-type: none"> <li>{ <i>statements</i> }</li> <li>-Executes the statements as an atomic operation.</li> </ul> </li> <li>• Stream I/O:               <ul style="list-style-type: none"> <li>-<math>x \leftarrow s</math> reads from a stream <math>s</math>, and saves the result in <math>x</math>.</li> <li>-<math>e \rightarrow s</math> writes the expression <math>e</math> to stream <math>s</math>.</li> </ul> </li> <li>• Attribute Access:               <ul style="list-style-type: none"> <li>-<math>x.y</math> accesses attribute <math>y</math> defined as part of variable <math>x</math>.</li> </ul> </li> </ul>
--	--

## 2.2 Operators

An operator can have zero or more parameters. These parameters can be simple functions, or initial values. We start by discussing the data processing operators.

**Map** is an operator, that *synchronously* reads from  $k$  incoming streams, applies a user-defined function  $f : X_1 \times X_2 \times \dots \times X_k \rightarrow Y$  on the read values, and writes the computed values to an outgoing stream. The operator is parametrized by the user-defined function.  $\text{MAP}(f) : \mathbf{S}\langle X_1 \rangle \times \dots \times \mathbf{S}\langle X_k \rangle \rightarrow \mathbf{S}\langle Y \rangle$

$$\mathbf{loop} : f(\leftarrow S_{\text{in}}^1, \dots, \leftarrow S_{\text{in}}^k) \rightarrow S_{\text{out}}$$

**Reduce** is an operator that has an internal state  $u : U$ . The operator reads from an incoming stream and applies a user defined function  $g : U \times X \rightarrow U \times Y$ . This function takes the saved state and the read value. Then, it computes a new state, and an output value. The operator updates the internal state, and writes the computed value to the output stream. The operator is parametrized by the function  $g$ , and an initial state  $u_0$ .  $\text{REDUCE}(g, u_0) : \mathbf{S}\langle X \rangle \rightarrow \mathbf{S}\langle Y \rangle$

$$\begin{aligned} \mathbf{state} \quad & u = u_0 \\ \mathbf{loop} : & u, y = g(u, \leftarrow S_{\text{in}}) \\ & y \rightarrow S_{\text{out}} \end{aligned}$$

**Copy** reads from an incoming stream, and synchronously duplicates the read value to all outgoing streams. **Copy** has no parameters.  $\text{COPY}() : \mathbf{S} \rightarrow \mathbf{S}^n$

$$\begin{aligned} \mathbf{loop} : & x \leftarrow S_{\text{in}} \\ & x \rightarrow S_{\text{out}}^i \quad \text{for all } i \leq n \end{aligned}$$

**Filter** has one incoming stream and two outgoing streams  $S_{\text{out}}^1$  and  $S_{\text{out}}^2$ . It applies a user-defined predicate  $\theta : X \rightarrow \text{boolean}$ , on the incoming values. It then writes values with  $\theta$  true to  $S_{\text{out}}^1$ , and values with  $\theta$  false to  $S_{\text{out}}^2$ .  $\text{FILTER}(\theta) : \mathbf{S} \rightarrow \mathbf{S}^2$

$$\begin{aligned} \text{loop} : x \leftarrow S_{\text{in}} \\ \text{if } \theta(x) \text{ then } x \rightarrow S_{\text{out}}^1 \text{ else } x \rightarrow S_{\text{out}}^2 \end{aligned}$$

**Ground** ends an incoming stream.  $\text{GROUND} : \mathbf{S} \rightarrow \emptyset$

$$\text{loop} : \leftarrow S_{\text{in}}$$

Now, we will discuss the rate controlling operators:

**Latch** has one incoming stream  $S_{\text{in}}$  and two outgoing streams  $S_{\text{out}}^1$  and  $S_{\text{out}}^2$ . It reads from  $S_{\text{in}}$ , synchronously writes to  $S_{\text{out}}^2$ , and asynchronously writes to  $S_{\text{out}}^1$ . It performs the asynchronous write by saving the most-recent incoming reading, and writing it to  $S_{\text{out}}^1$ , whenever it is possible. So  $S_{\text{in}}$  and  $S_{\text{out}}^1$  have different data rates.  $\text{LATCH}() : \mathbf{S} \rightarrow \mathbf{S}^2$ .

$$\begin{aligned} \text{state } u = \text{nil} \\ \text{loop}_1 : x \leftarrow S_{\text{in}} \quad \text{loop}_2 : \{u \rightarrow S_{\text{out}}^1\} \\ \{u = x ; x \rightarrow S_{\text{out}}^2\} \end{aligned}$$

**Cut** is similar to **Latch**, but it writes the incoming values only once to the asynchronous output stream. For the extra writes, in case  $S_{\text{out}}^1$  has a higher data rate, **nil** is used.  $\text{CUT}() : \mathbf{S} \rightarrow \mathbf{S}^2$

$$\begin{aligned} \text{state } u = \text{nil} \\ \text{loop} : x \leftarrow S_{\text{in}} \quad \text{loop} : \{y = u ; u = \text{nil}\} \\ \{u = x ; x \rightarrow S_{\text{out}}^2\} \quad y \rightarrow S_{\text{out}}^1 \end{aligned}$$

**Mult** has  $k$  incoming streams  $S_{\text{in}}^k$ , and one output stream  $S_{\text{out}}$ . The operator reads one value at a time from each incoming stream, forms a vector  $(x_1, \dots, x_k)$ , and synchronously writes this vector to the outgoing stream.  $\text{MULT}() : \mathbf{S}^k \rightarrow \mathbf{S}$

$$\text{loop} : \begin{bmatrix} \leftarrow S_{\text{in}}^1 \\ \dots \\ \leftarrow S_{\text{in}}^k \end{bmatrix} \rightarrow S_{\text{out}}$$

**Left-Mult** is similar to **Mult**; however, it has only two incoming streams  $S_{\text{in}}^1$  and  $S_{\text{in}}^2$ . It latches on  $S_{\text{in}}^2$  to make the outgoing data rate depends only on the rate of  $S_{\text{in}}^1$ , and independent of  $S_{\text{in}}^2$ .  $\text{LEFT-MULT} : \mathbf{S}^2 \rightarrow \mathbf{S}$

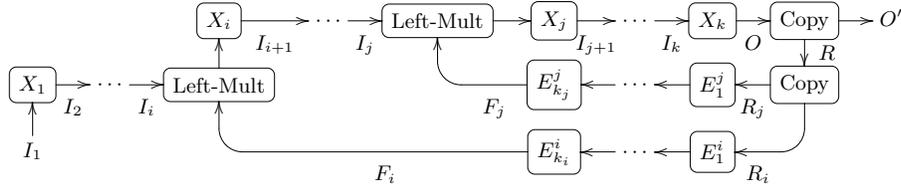
$$\begin{aligned} S^1, S^2 = \text{LATCH}(S_{\text{in}}^2) ; \text{GROUND}(S^2) \\ \text{loop} : \begin{bmatrix} \leftarrow S_{\text{in}}^1 \\ \leftarrow S^1 \end{bmatrix} \rightarrow S_{\text{out}} \end{aligned}$$

**Add** asynchronously merges together  $k$  incoming streams  $S_{\text{in}}^k$  into one outgoing stream.  $\text{ADD} : \mathbf{S}^k \rightarrow \mathbf{S}$ .

for  $j \leq k$   
 $\mathbf{loop}_j : x \leftarrow S_{\text{in}}^j$   
 $x \rightarrow S_{\text{out}}$

### 2.3 Feedback Control

Feedback control is an essential task in several online computer vision algorithms [14, 22–25] that perform parameter tuning, or iterative optimization. These algorithms evaluate the current output results to enhance the future outputs. In this section, we extend the stream algebra of [21] by providing a formal and abstract description of feedback control in computer vision pipelines.



**Fig. 1.** An example of multi-point feedback control. Arrows show stream flow directions. Letters on arrows represent stream names.  $I_1$ , and  $O'$  indicate input and output streams, respectively.

We assume an input stream  $I_1$ , a vision pipeline with a sequence of operators  $X_1, \dots, X_k$ , and an output stream  $O$ . Each operator  $X_j$  has an input stream  $I_j$ , and an output stream  $I_{j+1}$ . In order to define feedback control for operator  $X_j$ , we assume a return stream  $R_j$ , a sequence of evaluation operators  $E_1^j, \dots, E_{k_j}^j$ , and a feedback stream  $F_j$ . Given these assumptions; we describe the feedback loop in a vision pipeline as in figure 1. We discuss two types of feedback control, single-point and multi-point. In single-point, we only control one operator  $X_j$ . So, the pipeline has one return stream  $R$ , and one feedback stream  $F$ . In order to obtain the return stream  $R$ , we can apply either the COPY, FILTER, or CUT operator on the output stream  $O$ . This is represented by the equation,

$$R, O' \triangleq \text{COPY}()(O). \quad (1)$$

$O'$  now represents the output. We can next apply a sequence of evaluation operators  $E_1, \dots, E_n$  on  $R$  to compute the feedback stream  $F$ . If the feedback loop performs parameter tuning, then  $F$  represents a stream of new parameters. We can then apply a LEFT-MULT operator to attach the new parameters to the

input stream  $I_j$ , which updates the internal parameters of the pipeline operator  $X_j$ . This is represented by the equation,

$$I'_j \triangleq \text{LEFT-MULT}()(I_j, F). \quad (2)$$

$I'_j$  in this case represents the input to operator  $X_j$ . Note that, the feedback loop that we just defined is synchronized with the original vision pipeline. If we replace COPY in equation 1 by CUT, then we have an asynchronous feedback loop that evaluates samples from the output stream  $O$ . Furthermore, in case of iterative optimization that reprocesses the output, we can replace LEFT-MULT in equation 2 by an ADD operator. In this case, the stream  $I'_j$  will have interleaved elements from the streams  $I_j$  and  $F$ .

In multi-point feedback control, we can control more than one operator. This is performed by applying COPY on the return stream  $R$  several times to get  $m$  duplicate streams, which we use to define the feedback loops of  $m$  operators. Figure 1 shows multi-point feedback for two operators  $X_i$  and  $X_j$ . Each operator has a distinct set of evaluation operators to construct its feedback stream.

### 3 Feedback Control in Computer Vision Pipelines

There is a large interest in developing online computer vision algorithms that use feedback control to perform parameter tuning or iterative optimization tasks. These tasks allow an online algorithm to adapt itself continuously to different scene contexts, or iteratively improve output results over time. In this section, we discuss two state-of-the-art algorithms [24, 14] that process vision streams and apply feedback control to perform parameter tuning and iterative optimization. Without loss of generality, we will discuss how we can effectively express the feedback control of these algorithms using our algebraic extensions.

#### 3.1 Online Adaptation of Tracking Parameters

Online tracking of moving objects is one of the fundamental problems in computer vision, and several trackers have been proposed. However, if the input video stream has an unknown scene, it becomes difficult to select a tracking algorithm. Chau *et al.* [24] proposed a method for online parameter tuning that can adapt a tracking algorithm to scene changes. This method works in two phases: an offline training phase, and an online control phase.

In the offline training phase, the algorithm takes as input, a set of training videos together with annotated moving objects, annotated trajectories, and a tracking algorithm with parameters. The method tracks people as the moving objects using the appearance based tracking algorithm of [26], which is controlled by six parameters. For each training video, the algorithm extracts context features from every frame. The context features are a vector of six elements that describes the density of moving objects, their occlusion level, and appearance characteristics. The algorithm then segments the video into consecutive clips, based on similarity of the context feature vectors. For each clip, the

algorithm performs parameter optimization to select the best parameter values of the given tracking algorithm. Then, a clustering step is performed to cluster contexts from all training videos. Afterwards, the best tracking parameters are defined for each cluster. Finally, the context clusters together with their best parameters are stored in a database  $D$ .

In the online stage, the algorithm takes as input, a video stream  $V = \{V_i | i = 0, 1, 2, \dots\}$ . For every frame  $V_i \in V$ , people is detected using the HOG-based detection algorithm in [27]. This generates the objects stream  $B = \{B_i | i = 0, 1, 2, \dots\}$ , where every  $B_i$  represents a list of detected objects (people) in frame  $V_i$ . The algorithm then records the detected objects from every frame in a temporal window of interval  $\Delta t_1$ . After that, the method attaches with every frame  $V_i \in V$ , the recent list of recorded objects up to frame  $V_i$ . This generates the stream  $U = \{U_i | i = 0, 1, 2, \dots\}$ , where  $U_i = (V_i, B_i)$ . The method also defines the parameters feedback stream  $F = \{F_i | i = 0, 1, 2, \dots\}$ , where every  $F_i$  is a vector of tracking parameters. We will see later how this stream is generated. Every  $F_i$  is attached to  $U_i$  to generate the tracker input stream  $Q = \{Q_i | i = 0, 1, 2, \dots\}$ , where  $Q_i = (V_i, B_i, F_i)$ . The tracking algorithm takes the stream  $Q$  as input and updates the tracking parameters using  $F_i \in Q_i$ . It also generates a list of trajectories  $T_i$  for every frame. The tracker attaches the trajectories  $T_i$  to  $(V_i, B_i) \subset Q_i$ , and constructs the output trajectories stream  $J = \{J_i | i = 0, 1, 2, \dots\}$ , where  $J_i = (V_i, B_i, T_i)$ . The algorithm in [24], then defines a feedback control loop that takes the output stream  $J$  and records the video frames  $V_i \in J_i$  in a temporal window of interval  $\Delta t_2$ . This defines a clip stream  $C = \{C_i | i = 0, 1, 2, \dots\}$ . This stream is used with  $J$  to define the loopback stream  $H = \{H_i | i = 0, 1, 2, \dots\}$ , where  $H_i = (V_i, B_i, T_i, C_i)$ . Now, for every  $H_i$ , the algorithm uses  $(V_i, B_i, T_i) \subset H_i$  to calculate two error scores: An object interaction score  $s_1$  to calculate the overlap of objects, and a tracking error score  $s_2$  to measure the tracking quality. Given two thresholds  $Th_1$  and  $Th_2$ , if  $s_1 > Th_1$  and  $s_2 > Th_2$ , the algorithm declares an error and retrieves the best tracking parameters suitable to the context defined by the clip  $C_i \in H_i$ , from the database  $D$ . Otherwise, we continue using the current parameters. This generates the feedback stream  $F$  that was used previously together with the stream  $U$  to define the tracker input stream  $Q$ . This illustrates the feedback control loop of [24].

Now, we will describe the online control stage of [24] in the algebra of [21], and use our algebraic extensions to define the feedback control loop. We start by defining the following data types,

```

Frame : 2DImage;   Video : S ⟨Frame⟩;   Clip : LIST ⟨Frame⟩
Histogram : LIST ⟨R⟩;   Object : R8 × Histogram;   Params : R6
FrameInfo : Frame × LIST ⟨Object⟩;   Trajectory : LIST ⟨R2⟩;
TrackInput : Frame × LIST ⟨Object⟩ × Params
TrackInfo : Frame × LIST ⟨Object⟩ × LIST ⟨Trajectory⟩
LoopBack : Frame × LIST ⟨Object⟩ × LIST ⟨Trajectory⟩ × Clip

```

where a **Frame** is a 2D image, a **Video** is a stream of frames, a **Clip** is a list of frames, and a **Histogram** is a vector of values. An **Object** is a pair  $(a, b)$ , where  $b : \text{Histogram}$ , and  $a : \mathbb{R}^8$  is a vector that represents the following object features: 2D shape ratio, 2D area, color covariance in RGB, and dominant color in RGB. **Params** is a vector that represents the 6 parameters of the tracking algorithm [26]. **FrameInfo** is a pair  $(v, w)$ , where  $v : \text{Frame}$  and  $w : \text{LIST}(\text{Object})$ . A **TrackInput** is a 3 elements vector  $(v, w, p)$ , where  $p : \text{Params}$ . A **Trajectory** is a list of 2D points. **TrackInfo** is a 3 elements vector  $(v, w, e)$ , where  $e : \text{LIST}(\text{Trajectory})$ . Finally a **LoopBack** is a 4 elements vector  $(v, w, e, c)$ , where  $c : \text{Clip}$ .

Given an input video stream  $V \in \text{Video}$ , we copy  $V$  into two identical streams  $V_1, V_2$  using a **COPY** operator,

$$V_1, V_2 \triangleq \text{COPY}()(V). \quad (3)$$

We will now process  $V_1$  and return later to discuss the use of  $V_2$ . We define a function  $f_1 : \text{Frame} \rightarrow \text{LIST}(\text{Object})$  that detects objects from every frame  $V_i \in V_1$  using the HOG-based detection algorithm of [27]. This function can be used with a **MAP** operator to define the objects stream  $B : \mathbf{S}(\text{LIST}(\text{Object}))$ ,

$$B \triangleq \text{MAP}(f_1)(V_1). \quad (4)$$

We then define the function,

$$\begin{aligned} g_1 : \text{LIST}(\text{Object}) \times \text{LIST}(\text{Object}) &\rightarrow \text{LIST}(\text{Object}) \times \text{LIST}(\text{Object}) \\ g_1(u, x) = \{ &\text{for all } z \in u \\ &\text{if } (\text{now}() - \text{arrival-time}(z) \geq \Delta t_1) \text{ then} \\ &\quad u = u \ominus z \quad // \text{remove } z \text{ from } u \\ &u = u \oplus x \quad // \text{append } x \text{ to } u \\ &\text{return}(u, u) \quad \} \end{aligned}$$

This function maintains a window  $u$  of the most recent objects, in a time interval  $\Delta t_1$ . The function starts by deleting old objects. Then, it adds the new objects  $x$  to  $u$ , and returns the updated window as output. We can use the  $g_1$  function together with a **REDUCE** operator to generate the objects summary stream  $M : \mathbf{S}(\text{LIST}(\text{Object}))$ ,

$$M \triangleq \text{REDUCE}(g_1, \text{Empty-List})(B). \quad (5)$$

Now, we will go back and use the stream  $V_2$ . Note that this stream is a copy of the input video stream  $V$ . We synchronize the stream  $M$  with the stream  $V_2$  using a **MULT** operator to generate the stream  $U : \mathbf{S}(\text{FrameInfo})$ ,

$$U \triangleq \text{MULT}()(V_2, M). \quad (6)$$

At this step, we define the feedback stream  $F : \mathbf{S}(\text{Params})$ . We will show later how we generate this stream when we discuss feedback control. We synchronize

the stream  $F$  with the stream  $U$  using a LEFT-MULT operator to generate the stream  $Q : S \langle \text{TrackInput} \rangle$ ,

$$Q \triangleq \text{LEFT-MULT}()(U, F). \quad (7)$$

We then define the tracking function  $f_2 : \text{TrackInput} \rightarrow \text{TrackInfo}$ . This function takes the vector  $(v, w, p) : \text{TrackInput}$ , and applies the tracking algorithm on  $v$  and  $w$  using the given parameters  $p$ . The function then outputs the object trajectories which we attach to the pair  $(v, w)$ , to generate the output  $y : \text{TrackInfo}$ . We can use this function together with a MAP operator to process the  $Q$  stream and generate the trajectories stream  $J : S \langle \text{TrackInfo} \rangle$ ,

$$J \triangleq \text{MAP}(f_2)(Q). \quad (8)$$

Now, we discuss how we use our algebraic description to express the feedback control of [24]. We start by evaluating the equation,  $R, J' \triangleq \text{COPY}()(J)$ , to copy  $J$  into the output stream  $J'$  and the return stream  $R$ . Then, we define the function,

$$g_2 : \text{Clip} \times \text{TrackInfo} \rightarrow \text{Clip} \times \text{LoopBack}$$

$$g_2(u, x) = \{ \text{for all } z \in u \\ \text{if } (\text{now}() - \text{arrival-time}(z) \geq \Delta t_2) \text{ then} \\ \quad u = u \ominus z \quad // \text{remove } z \text{ from } u \\ u = u \oplus x.v \quad // \text{append frame } x.v \text{ to } u \\ y = x \oplus u \quad // \text{append clip } u \text{ to } x \\ \text{return}(u, y) \}$$

The function  $g_2$  maintains a clip  $u$  that has the most recent frames in a time interval  $\Delta t_2$ . It is similar to the function  $g_1$ , however  $g_2$  appends the recorded clip to the input  $x$ , to form the output  $y : \text{LoopBack}$ . We use this function with a REDUCE operator to process the return stream  $R$ , and generate the loopback stream  $H : S \langle \text{LoopBack} \rangle$ ,

$$H \triangleq \text{REDUCE}(g_2, \text{Empty-List})(R). \quad (9)$$

The algorithm in [24] defines two scoring functions to calculate the tracking errors. These functions are  $f_3 : \text{FrameInfo} \rightarrow \mathbb{R}$  for object interaction score, and  $f_4 : \text{TrackInfo} \rightarrow \mathbb{R}$  for tracking error score. We use these functions to define another function,

$$g_3 : \text{Params} \times \text{LoopBack} \rightarrow \text{Params} \times \text{Params}$$

$$g_3(u, x) = \{ s_1 = f_3(x.v, x.w) \\ s_2 = f_4(x.v, x.w, x.e) \\ \text{if } (s_1 > Th_1 \text{ and } s_2 > Th_2) \text{ then} \\ \quad p = \text{search-db}(x.c, D) \\ \quad \text{return}(p, p) \\ \text{return}(u, u) \}$$

Note that this function starts by calculating the error scores of the current frame. If the scores  $s_1$  and  $s_2$  are larger than the given thresholds, we search the database  $D$  for the best parameters that meet the context of the most recent clip  $x.c$ , and output the new parameters. Otherwise, we return the old parameters  $u$ . This function can be used together with a REDUCE operator to process the loopback stream  $H$ , and generate the feedback stream  $F : S \langle \text{Params} \rangle$ ,

$$F \triangleq \text{REDUCE}(g_3, \text{Initial-Params})(H). \quad (10)$$

Remember that we used the stream  $F$  together with the stream  $U$  as input to the LEFT-MULT operator in equation 7 to define the tracker input stream  $Q$ . In this way, we use our definition of single-point feedback control to continuously update the tracker parameters.

### 3.2 Iterative Optimization for Aligning Photo Streams

Recently, there has been a large interest in the analysis of Web photo streams. This interest is driven by the existence of several photo hosting websites that contain huge amounts of personal photo collections. Kim *et al.* [14] proposed a recent algorithm to build common storylines from Flickr’s photo streams. Their algorithm takes as input, a set of  $n$  photo streams  $I = \{I_k | k = 1, 2, 3..n\}$  from different Flickr users that share a common user activity. Each photo in a stream  $I_k$  stores its capture time, a spatial pyramid histogram as a visual descriptor, and a set of foreground regions that is initially empty. Every stream  $I_k$  is divided into a sequence of photo blocks  $B_k = \{B_{\{k,i\}} | i = 0, 1, 2.. \}$ , where each block  $B_{\{k,i\}}$  represents the photos taken by a user over a certain period of time  $\Delta t_1$  (for example, a day). Each block also stores the earliest and latest capture times of its photos. The algorithm then iterates between two tasks, an alignment task, and a cosegmentation task.

In the alignment task, the algorithm starts by reading one block from each stream, and constructs the block-list stream  $L = \{L_i | i = 0, 1, 2.. \}$ . Every  $L_i \in L$  is a list of blocks. Then, the algorithm selects from every list  $L_i$ , a set of blocks  $E_i$  that overlap in a timeline by at least a period  $\Delta t_2$  (for example, an hour). We also attach to  $E_i$ , an iteration number  $N_i$ , that is initially zero. This generates the overlapped blocks stream  $Q = \{Q_i | i = 0, 1, 2.. \}$ , where  $Q_i = (E_i, N_i)$ . The algorithm defines a feedback stream  $F$  which is of the same type as  $Q$ . We will see later how this stream is constructed. We add the elements of the two streams  $Q$  and  $F$  together to construct the interleaved stream  $P = \{P_j | j = 0, 1, 2.. \}$ , where  $P_j = (E_j, N_j)$ . Note that the stream  $P$  contains interleaved elements from the two streams  $Q$  and  $H$ . For each block  $b$  of  $E_j \in P_j$ , the algorithm calculates the similarity of  $b$  to other blocks in  $E_j$ . This is performed for two blocks  $(b_1, b_2) \subset E_j$  by first finding for each photo  $x \in b_1$ , its best visually similar neighbor  $y \in b_2$ . This similarity is calculated using a distance function  $f_4 : \text{Photo} \times \text{Photo} \rightarrow \mathbb{R}$ . If both  $x$  and  $y$  have foreground regions defined,  $f_4$  returns the distance between the histograms of these regions. Otherwise,  $f_4$  returns the distance between the spatial pyramid histograms of photos. The algorithm also defines another

distance function  $f_5 : \text{Time} \times \text{Time} \rightarrow \mathbb{R}$ , that calculates the difference between the capture times of two photos. The algorithm then measures the similarity between two blocks  $b_1$  and  $b_2$  using an energy function  $f_6 : \text{Block} \times \text{Block} \rightarrow \mathbb{R}$ , that sums the similarity defined by  $f_4$  and  $f_5$  along the correspondent photos between  $b_1$  and  $b_2$ . The method then maps every list of blocks  $E_j \in P_j$  into a graph  $G_j : \text{Graph} \langle \text{Block}, \text{Block} \times \text{Block} \rangle$ . This graph has the blocks of  $E_j$  as vertices. An edge exists between two blocks, if they are best similar to each other, in other words, they have the minimal distance to each other according to  $f_6$ . This generates the blocks-graph stream  $Z = \{Z_j | j = 0, 1, 2, \dots\}$ , where  $Z_j = (G_j, N_j)$ . This stream defines the output of the alignment step.

The cosegmentation step takes as input, the blocks-graph stream  $Z$ . This step maps every graph  $G_j \in Z_j$  into a new graph  $Y_j : \text{Graph} \langle \text{Photo}, \text{Photo} \times \text{Photo} \rangle$ . This is performed by first collecting all photos from the blocks of  $G_j$ . These photos define the vertices of  $Y_j$ . Then, the algorithm adds an edge between two photos  $(x, y)$  if they are a correspondent pair of two different blocks. In addition, for every photo  $x$  in a block  $b$ , the method adds edges to its  $k$  nearest neighbors from the same block  $b$ . We append the graph  $Y_j$  to  $Z_j$  to construct the photos-graph stream  $M = \{M_j | j = 0, 1, 2, \dots\}$ , where  $M_j = (G_j, N_j, Y_j)$ . We also increment  $N_j$  by 1. Afterwards, the algorithm in [28] defines for each photo in the vertices of  $Y_j \in M_j$ , a set of  $m$  foreground regions. This is performed by applying the cosegmentation technique of [28] and belief propagation between each photo and its neighbors in  $Y_j$ . If a photo already have foreground regions defined, then they are enhanced.

The algorithm of [14] iterates between the alignment task and the cosegmentation task. This is achieved by defining a feedback loop that first filters the photos-graph stream  $M$  based on the iteration number  $N_j \in M_j$ . If  $N_j \geq N_{stop}$  then  $M_j$  is sent to the output, Otherwise, it is sent to the feedback loop.  $N_{stop}$  defines the maximum number of iterations. This defines two streams  $M'$  and  $R$  with similar type to stream  $M$ , where  $M'$  is the output stream, and  $R$  is the return stream. The stream  $R$  is then mapped to the feedback stream  $F = \{F_l | l = 0, 1, 2, \dots\}$ , where  $F_l = (E_l, N_l)$ .  $E_l$  is the list of blocks in the vertices of  $G_l \in R_l$  and  $N_l$  is the iteration number. Remember that the feedback stream  $F$  was added to the stream  $Q$  to define the interleaved stream  $P$  in the alignment step. Note also that the block photos in  $F$  have foreground regions defined. These regions enhance the matching of the feedback blocks in the alignment step. Consequently, This improves the output of the cosegmentation step, which closes the feedback loop defined by [14] for iterative optimization.

Now, we will express the iterative optimization of [24] using our algebraic extensions for feedback control. We define the following data types,

```
Shape : LIST  $\langle \mathbb{R}^2 \rangle$ ;   Region : Shape  $\times$  Histogram
Photo : 2DImage  $\times$  Time  $\times$  Histogram  $\times$  LIST (Region)  $\times$   $\mathbb{R}$ 
Block : LIST (Photo)  $\times$  Time2;   BlocksInfo : LIST (Block)  $\times$   $\mathbb{R}$ ;
BlocksGraphInfo : GRAPH (Block, Block  $\times$  Block)  $\times$   $\mathbb{R}$ 
```

**PhotosGraphInfo** : **BlocksGraphInfo**  $\times$  **GRAPH**  $\langle$ Photo, Photo  $\times$  Photo $\rangle$

Where, a **Shape** is a list of 2D points. A **Region** is a 2D vector  $(s, h)$ , where  $s$  : **Shape**, and  $h$  : **Histogram** is the region descriptor. We reuse the definition of **Histogram** from the previous section. A **Photo** is a vector of five variables  $(a, t, h, r, nn)$ , where  $a$  : **2DImage**,  $t$  : **Time** is the capture time,  $h$  : **Histogram** is a visual descriptor,  $r$  : **LIST**  $\langle$ **Region** $\rangle$  is a list of foreground regions (initially empty), and  $nn$  is the index of the nearest neighbor photo. A **Block** is a 3D vector  $(b, t_1, t_2)$ , where  $b$  : **LIST**  $\langle$ **Photo** $\rangle$ ,  $t_1$  : **Time** is the earliest capture time of photos in  $b$ , and  $t_2$  is the latest capture time. A **BlocksInfo** is a vector  $(w, itr)$ , where  $w$  : **LIST**  $\langle$ **Block** $\rangle$ , and  $itr$  :  $\mathbb{R}$  indicates the iteration number. A **BlocksGraphInfo** is a vector  $(c_1, itr)$ , where  $c_1$  is a graph on a set of blocks. Finally, **PhotosGraphInfo** is a 2D vector  $(q, c_2)$ , where  $q$  : **BlocksGraphInfo**, and  $c_2$  is a graph on a set of photos.

For simplicity, we will consider that we have three input photo streams  $I = \{I_k | k = 1..3\}$ . We now define the following function,

$g_4 : \mathbf{Block} \times \mathbf{Photo} \rightarrow \mathbf{Block} \times \mathbf{Block}$

$$g_4(u, x) = \{ \text{if duration}(u) \geq \Delta t_1 \text{ then} \\ \quad u' = \emptyset; y = u \\ \text{else} \\ \quad u' = u \oplus x \quad // \text{append } x \text{ to Block } u \\ \quad y = \emptyset \\ \text{return}(u', y) \}$$

This function keeps appending the incoming photos to a block  $u'$  and setting the output  $y$  to an empty block. When the duration of the block  $\Delta t = u'.t_2 - u'.t_1$  is larger than a certain interval  $\Delta t_1$  (A day in [14]), the function copies the block  $u'$  to the output  $y$  and resets  $u'$  back to an empty block. We can use this function together with a **REDUCE** and **FILTER** operators to process every stream  $I_k \in I$  and generate a corresponding stream  $B_k : S \langle \mathbf{Block} \rangle$ ,

$$B_k, E_k \triangleq \mathbf{FILTER}(\lambda x : |x| \neq 0) \circ \mathbf{REDUCE}(g_4, \mathbf{Empty-List})(I_k). \quad (11)$$

$$\mathbf{GROUND}()(E_k). \quad (12)$$

Where the  $\circ$  operator is defined by [21] as a composition operator, that supplies the output stream from the right operand as an input stream to the left operand. The **FILTER** operator removes the empty blocks from the output stream of the **REDUCE** operator. These empty blocks define the stream  $E_k$  which is ignored by a **GROUND** operator. We can now synchronize the three streams  $B_1, B_2$  and  $B_3$  using a **MULT** operator to construct the block-list stream  $L : S \langle \mathbf{LIST} \langle \mathbf{Block} \rangle \rangle$ ,

$$L \triangleq \mathbf{MULT}()(B_1, B_2, B_3). \quad (13)$$

We then define the function  $f_7 : \text{LIST} \langle \text{Block} \rangle \rightarrow \text{BlocksInfo}$ . This function selects from every list  $L_i \in L$ , blocks that overlap in a timeline by at least a period  $\Delta t_2$  (An hour in [14]). The function also attaches to the list of selected blocks, an iteration number, which is initially zero. We can use  $f_7$  together with a MAP operator to define the stream  $Q : S \langle \text{BlocksInfo} \rangle$ ,

$$Q \triangleq \text{MAP}(f_7)(L). \quad (14)$$

At this step, we define the feedback stream  $F : S \langle \text{BlocksInfo} \rangle$ . We will show later how we generate this stream when we discuss the feedback loop. We add the two streams  $Q$  and  $F$  using an ADD operator to generate the stream  $P : S \langle \text{LIST} \langle \text{Block} \rangle \rangle$ ,

$$P \triangleq \text{ADD}()(Q, F). \quad (15)$$

Now, we define two functions  $f_8 : \text{BlocksInfo} \rightarrow \text{BlocksGraphInfo}$ , and  $f_9 : \text{BlocksGraphInfo} \rightarrow \text{PhotosGraphInfo}$ . The function  $f_8$  maps a blocks list in **BlocksInfo** into a graph of blocks, where an edge exists between two blocks if they are nearest neighbors according to the function  $f_6$ . Note that  $f_6$  gives better similarity distance if block photos have good foreground regions. The function  $f_9$ , on the other hand, maps a graph of blocks into a graph of photos, and increments the iteration number by 1. Cosegmentation is applied between each photo and its neighbors in the photos graph to define foreground regions or enhance existing ones. We can use the two functions  $f_8$  and  $f_9$  together with two MAP operators to define the photos-graph stream  $M : S \langle \text{PhotosGraphInfo} \rangle$ ,

$$M \triangleq \text{MAP}(f_9) \circ \text{MAP}(f_8)(P). \quad (16)$$

Now, we will discuss the feedback control of [14] to perform iterative optimization using our algebraic extensions. We start by applying a FILTER operator on the stream  $M$  to define the output stream  $M'$  and the return stream  $R$ ,

$$M', R \triangleq \text{FILTER}(\lambda x : x.q.itr \geq N_{stop})(M). \quad (17)$$

We access the attribute  $q : \text{BlocksGraphInfo}$  from each element  $M_j \in M$ . Then, we test if the iteration number of  $q$  reached the maximum number of iterations. If the test is true  $M_j$  is sent to  $M'$ , if not, it is sent to  $R$ . We then define a function  $f_{10} : \text{PhotosGraphInfo} \rightarrow \text{BlocksInfo}$  which maps every element in  $R$  into an element of type **BlocksInfo**. This is performed by converting the vertices of the blocks-graph attached to **PhotosGraphInfo** into a list of blocks, and copying the current iteration number. The function  $f_{10}$  together with a MAP operator define the feedback stream  $F : S \langle \text{BlocksInfo} \rangle$ ,

$$F \triangleq \text{MAP}(f_{10})(M_2). \quad (18)$$

Remember that we used the stream  $F$  together with the stream  $Q$  as input to the ADD operator in equation 15 to define the stream  $P$ . This shows how our definition of single-point feedback control describes iterative optimization in [14].

## 4 Discussions

The examples discussed in the previous section, demonstrate that we can effectively express feedback control in online computer vision algorithms using our abstract algebraic description presented in section 2.3. In each discussed algorithm, we first express its vision pipeline by a set of algebraic equations based on the operators defined by [21]. Then, we use our abstract definition of single-point feedback control to describe both the parameter tuning and the iterative optimization tasks. For example, the algebraic description of the online tracking example [24] shows a flexible integration of our single-point feedback to tune tracking parameters. Additionally, We can easily extend this single-point feedback control to multi-point feedback control (Figure 1). We can use this, for example, to control both the parameters of tracker and the parameters of the HOG-based people detector algorithm [27] used by [24] to detect moving objects. So, our multi-point description can scale up feedback to control several operators of a vision pipeline. Furthermore, if the vision pipeline has several output streams, we may extend our multi-point feedback to process multiple return streams. In this case, we can use rate control operators such as CUT and LATCH to replace COPY in figure 1, and define asynchronous feedback control that does not affect the flow rate of the feedforward pipeline. This also motivates that we may perform multi-point feedback between different computer vision pipelines, that work at separate data rates, without affecting one another.

The discussed iterative optimization example [14] suggests that our formal feedback description may be useful for other tasks such as adaptive learning and incremental evaluation. In addition, our feedback description may be used for implementing blocking resolution in vision pipelines. For example, we can make every pipeline operator attach its estimated runtime to the output stream. Then, these runtimes can be collectively monitored by feedback control to send back appropriate actions to the operators. So, our formal feedback description opens new directions for tasks such as real-time debugging, performance monitoring and bottleneck identification.

## 5 Conclusion

This paper develops abstract and formal methods for describing and implementing feedback control in computer vision pipelines—online vision algorithms that process images and videos (vision streams). These formal methods build upon an existing stream algebra to flexibly integrate feedback control to vision pipelines. We have demonstrated our formal methods for two state-of-the-art online computer vision algorithms that implement feedback control for parameter tuning and iterative optimization. We also discussed how our methods can scale up feedback to control different stream operators of the vision pipeline. Our work opens up new research directions to study real-time debugging, dynamic reconfiguration, bottleneck identification, adaptive learning, and performance tuning of vision pipelines.

## References

1. Zhao, B., Fei-Fei, L., Xing, E.: Online detection of unusual events in videos via dynamic sparse coding. In: CVPR, Colorado Springs (2011) 3313–3320
2. Helala, M., Pu, K., Qureshi, F.: Road boundary detection in challenging scenarios. In: AVSS. (2012) 428–433
3. Meghdadi, A., Irani, P.: Interactive exploration of surveillance video through action shot summarization and trajectory visualization. *IEEE Trans. on Visualization and Computer Graphics* **19** (2013) 2119–2128
4. Yenikaya, S., Yenikaya, G., Düven, E.: Keeping the vehicle on the road: A survey on on-road lane detection systems. *ACM Comput. Surv.* **46** (2013) 2:1–2:43
5. Ozcanli, O., Dong, Y., Mundy, J., Webb, H., Hammoud, R., Victor, T.: Automatic geo-location correction of satellite imagery. In: IEEE CVPR Workshops. (2014)
6. Wischounig-Strucl, D., Quartisch, M., Rinner, B.: Prioritized data transmission in airborne camera networks for wide area surveillance and image mosaicking. In: IEEE CVPR Workshops. (2011) 17–24
7. Yuping, L., Medioni, G.: Map-enhanced uav image sequence registration and synchronization of multiple image sequences. In: IEEE CVPR, Minneapolis, Minnesota, USA (2007) 1–7
8. Ryoo, M.S.: Human activity prediction: Early recognition of ongoing activities from streaming videos. In: ICCV, Barcelona, Spain (2011) 1036–1043
9. Lu, C., Shi, J., Jia, J.: Online robust dictionary learning. In: IEEE CVPR. (2013) 415–422
10. Xuand, C., Xiong, C., Corso, J.: Streaming hierarchical video segmentation. In: ECCV. Volume VI. (2012) 626–639
11. Loy, C., Hospedales, T., Xiang, T., Gong, S.: Stream-based joint exploration-exploitation active learning. In: CVPR. (2012) 1560–1567
12. Harbi, N., Gotoh, Y.: Spatio-temporal human body segmentation from video stream. In: *Computer Analysis of Images and Patterns*. Volume 8047. Springer (2013) 78–85
13. Yang, J., Luo, J., Yu, J., Huang, T.: Photo stream alignment and summarization for collaborative photo collection and sharing. *IEEE Trans. on Multimedia* **14** (2012) 1642–1651
14. Kim, G., Xing, E.: Jointly aligning and segmenting multiple web photo streams for the inference of collective photo storylines. In: CVPR. (2013) 620–627
15. Chkodrov, G., Ringseth, P., Tarnavski, T., Shen, A., Barga, R., Goldstein, J.: Implementation of stream algebra over class instances, Google patents (2013)
16. Broy, M., Stefanescu, G.: The algebra of stream processing functions. *Theoretical Computer Science* **258** (2001) 99 – 129
17. Carlson, J., Lisper, B.: An event detection algebra for reactive systems. In: *Proceedings of the 4th ACM International Conference on Embedded Software*. (2004) 147–154
18. Demers, A., Gehrke, J., Hong, M., Riedewald, M., White, W.: A general algebra and implementation for monitoring event streams. Technical report, Cornell University (2005)
19. Shen, C., Little, J., Fels, S.: Towards OpenVL: Improving real-time performance of computer vision applications. In: *Embedded Computer Vision. Advances in Pattern Recognition*. Springer London (2009) 195–216
20. GStreamer. <http://gstreamer.freedesktop.org> (2014) Accessed: 2014-01-26.
21. Helala, M.A., Pu, K.Q., Qureshi, F.Z.: A stream algebra for computer vision pipelines. In: IEEE CVPR Workshops. (2014)
22. Kisilev, P., Freedman, D.: Parameter tuning by pairwise preferences. In: *BMVC*. (2010)
23. Sherrah, J.: Learning to adapt: A method for automatic tuning of algorithm parameters. In: *ACIVS*. (2010) 414–425

24. Chau, D., Badie, J., Bremond, F., Thomnat, M.: Online tracking parameter adaptation based on evaluation. In: IEEE International Conference on AVSS. (2013) 189–194
25. III, J.S., Ramanan, D.: Self-paced learning for long-term tracking. In: CVPR, Washington, DC, USA (2013) 2379–2386
26. Chau, D., Bremond, F., Thomnat, M.: A multi-feature tracking algorithm enabling adaptation to context variations. In: ICDP 2011. (2011) 1–6
27. Corvee, E., Bremond, F.: Body parts detection for people tracking using trees of histogram of oriented gradient descriptors. In: IEEE International Conference on AVSS. (2010) 469–475
28. Kim, G., Xing, E.: On multiple foreground cosegmentation. In: IEEE CVPR. (2012) 837–844