# REGULAR PROGRAMMING FOR QUANTITATIVE PROPERTIES OF DATA STREAMS

Rajeev Alur Dana Fisman Mukund Raghothaman ESOP 2016

University of Pennsylvania

· P2	· P2	· EOM	· P2
· EOD	· P1	· P2	· EOD
· P2	· EOD	· P2	· P2
· P1	· P2	· P2	· P2
· P1	· P2	· P1	· P2
· EOM	· P1	· P2	· EON
· P2	· P2	· EOM	· P1
· P2	· P1	· P1	·

 $\cdot$  What is the maximum number of daily requests for P2?

*iter-max(day-count*<sub>P2</sub>)

• What is the maximum number of daily requests for P2?

*iter-max(day-count*<sub>P2</sub>)

 $\cdot\,$  What is the number of requests for P1 in an average month?

*iter-avg*(*month-count*<sub>P1</sub>)

• What is the maximum number of daily requests for P2?

*iter-max(day-count*<sub>P2</sub>)

 $\cdot\,$  What is the number of requests for P1 in an average month?

*iter-avg(month-count*<sub>P1</sub>)

• What is the total number of requests for P1?

 $count_{P1} = iter-sum(day-count_{P1})$ 

• What is the maximum number of daily requests for P2?

*iter-max(day-count*<sub>P2</sub>)

 $\cdot\,$  What is the number of requests for P1 in an average month?

iter-avg(month-count<sub>P1</sub>)

• What is the total number of requests for P1?

 $count_{P1} = iter-sum(day-count_{P1})$ 

• What is the maximum of the total number of requests for P1 and P2?

 $max(count_{P1}, count_{P2})$ 

### Languages $\Sigma^* \to \text{bool} \equiv \text{Regular expressions}$ Cost functions $\Sigma^* \to \mathbb{R} \equiv \text{QREs}$

#### FUNCTION COMBINATORS

#### $a \mapsto d$

- · If w = a, then output d, otherwise undefined
- $\cdot \ \mbox{P1} \mapsto 0 \mbox{, P1} \mapsto 1 \mbox{, EOD} \mapsto 5 \mbox{, ...}$
- · Analogue of basic regular expressions

f else g

- · If f(w) is defined, then output f(w), otherwise output g(w)
- $\cdot$  Analogue of regular expression union

- f + g: Map w to f(w) + g(w)
- f g: Map w to f(w) g(w)
- $\max(f,g)$ : Map w to  $\max(f(w),g(w))$
- • •
- $\cdot op(f_1, f_2, \dots, f_k)$ : Map w to  $op(f_1(w), f_2(w), \dots, f_k(w))$
- $\cdot$  Analogue of regular expression intersection





split-op(f,g)

# split-plus(f, g)



split-op(f,g)

## split-plus(f,g), split-max(f,g)



split-op(f, g)

## split-plus(f,g), split-max(f,g), ..., split-op(f,g)





ATTEMPT 1

Basic functions:  $a \mapsto d$ Conditional choice: f else gConcatenation: split-op(f,g)Function iteration: iter-op(f)Cost operations:  $op(f_1, f_2, \dots, f_k)$  Basic functions:  $a \mapsto d$ Conditional choice: f else gConcatenation: split-op(f,g)Function iteration: iter-op(f)Cost operations:  $op(f_1, f_2, ..., f_k)$ 

• What about non-binary, non-commutative or non-associative operators?

Basic functions:  $a \mapsto d$ Conditional choice: f else gConcatenation: split-op(f,g)Function iteration: iter-op(f)Cost operations:  $op(f_1, f_2, ..., f_k)$ 

- What about non-binary, non-commutative or non-associative operators?
- $\cdot$  Are these operators "sufficient"? If we add more operators?

$$|-----]$$







• Attempt 1 annotates the parse tree with cost operations to obtain the computation tree



- Attempt 1 annotates the parse tree with cost operations to obtain the computation tree
- · What if QREs map input streams to computation trees?

Key Insight

QREs map input streams to terms over the cost domain

- · Terms contain parameters
- f(aaab) = 5 + p, where p is a parameter

Key Insight

QREs map input streams to terms over the cost domain

- · Terms contain parameters
- f(aaab) = 5 + p, where p is a parameter
- · Parameter substitution is also an operation

## Parameter subtitution, f[p/g]

- Say  $f(w) = t_f$ ,
- · and  $g(w) = t_g$
- $\cdot f[p/g](w) = t_f[p/t_g]$

Output

 $split(f \rightarrow^p g)$ 



 $split(f \rightarrow^p g)$ 



### **Simulating** *split-plus*(*f*, *g*)



· split-plus(f, g) = split( $f \rightarrow^{p} g'$ )

### **Simulating** *split-plus*(*f*, *g*)



- · split-plus(f, g) = split( $f \rightarrow^{p} g'$ )
- $\cdot g'(w) = p + g(w)$





Basic functions:  $a \mapsto d$ , regex  $\mapsto$  term Conditional choice: f else gConcatenation:  $split(f \rightarrow^p g)$ ,  $split(f \leftarrow^p g)$ Function iteration:  $iter \rightarrow(f)$ ,  $iter \leftarrow(f)$ Cost operations:  $op(f_1, f_2, \dots, f_k)$ , f[p/g] Basic functions:  $a \mapsto d$ ,  $regex \mapsto term$ Conditional choice: f else gConcatenation:  $split(f \rightarrow^p g)$ ,  $split(f \leftarrow^p g)$ Function iteration:  $iter \rightarrow(f)$ ,  $iter \leftarrow(f)$ Cost operations:  $op(f_1, f_2, \dots, f_k)$ , f[p/g]

- · QREs map string to terms
- $\cdot\,$  Structural operators decoupled from cost operators

- Was our choice of combinators ad-hoc? What functions can the formalism express?
- Given f and an input stream w, can we efficiently compute f(w)?

#### EXPRESSIVENESS

### Languages $\Sigma^* \to \text{bool} \equiv \text{Finite automata}$ Cost functions $\Sigma^* \to \mathbb{R} \equiv ?$

Regular languages have many appealing properties: many natural equi-expressive characterizations, robust closure properties, decidable analysis problems, practical utility

## COST REGISTER AUTOMATA [LICS 2013]



- · Finite state space, finitely many registers
- · Registers hold terms: Parameterized by set of operations
- · Equivalent to MSO definable string-to-term transducers
- · Closed under regular lookahead, input reversal, etc

#### Theorem

#### QREs can express exactly the same functions as CRAs

#### Taste of the completeness proof

- Piggy-back on DFA to regular expression translation Construct  $R^i(q,q')$ : all strings from q to q' while only traversing states less than  $q_i$
- Construct  $f^{i}(q, q', x)$ : express final value of register x as a DReX expression

#### Taste of the completeness proof

- · Capture data flows using "shapes"
- Construct a partial order over shapes, and use as basis for induction



#### FAST EVALUATION ALGORITHMS

#### Given a QRE f and an input stream w, find f(w)

- $\cdot\,$  QREs separate intent from evaluation
- If QREs are unambiguous, then f(w) can be computed with a single pass over w in time O(|w| · poly(|f|))







5 potential parse trees needed at each step

5 potential parse trees needed at each step

Idea 1: Statically bound number of potential parse trees:
O(poly(|f|)), irrespective of |w|









- Term compression ensures intermediate terms of bounded size
- Idea 2: If only operators are \*, +, -, min, max, *avg*, then space usage is polynomially bounded too!

#### CONCLUSION

- · Introduced Quantitative Regular Expressions (QRE)
- · Idea of function combinators
  - · Function descriptions are modular
  - $\cdot\,$  Regular parsing of the input data stream
- Simple, expressive programming model for stream processing, with strong theoretical foundations and fast evaluation algorithms

### Languages $\Sigma^* \to \text{bool} \equiv \text{Regular expressions}$ Cost functions $\Sigma^* \to \mathbb{R} \equiv \text{QREs}$

- $\cdot \:$  Also works for  $\Sigma^* \to \mathbb{N}$ ,  $\Sigma^* \to \mathbb{Q}$ , ...
- $\cdot \ \Sigma^* \to \mathbb{D}$  , for arbitrary cost domain  $\mathbb{D}$
- · Key insight: Generalizing to string-to-term transformations

- Expressively equivalent to regular cost functions / cost register automata
- Fast one-pass evaluation algorithms for unambiguous expressions
- Low space usage if only operations used are \*, +, -, min, max, avg

- Approximate evaluation algorithms for certain operations such as *iter-median* (Jointly with Sanjeev Khanna and Kostas Mamouras)
- $\cdot\,$  Exploring connections to streaming databases

## FIN!

QUESTIONS, COMMENTS, BRICKBATS?