# ASPEN: A Scalable In-SRAM Architecture for Pushdown Automata

**Kevin Angstadt**\*, Arun Subramaniyan\*, Elaheh Sadredini[†], Reza Rahimi[†], Kevin Skadron[†], Westley Weimer\*, Reetuparna Das\*

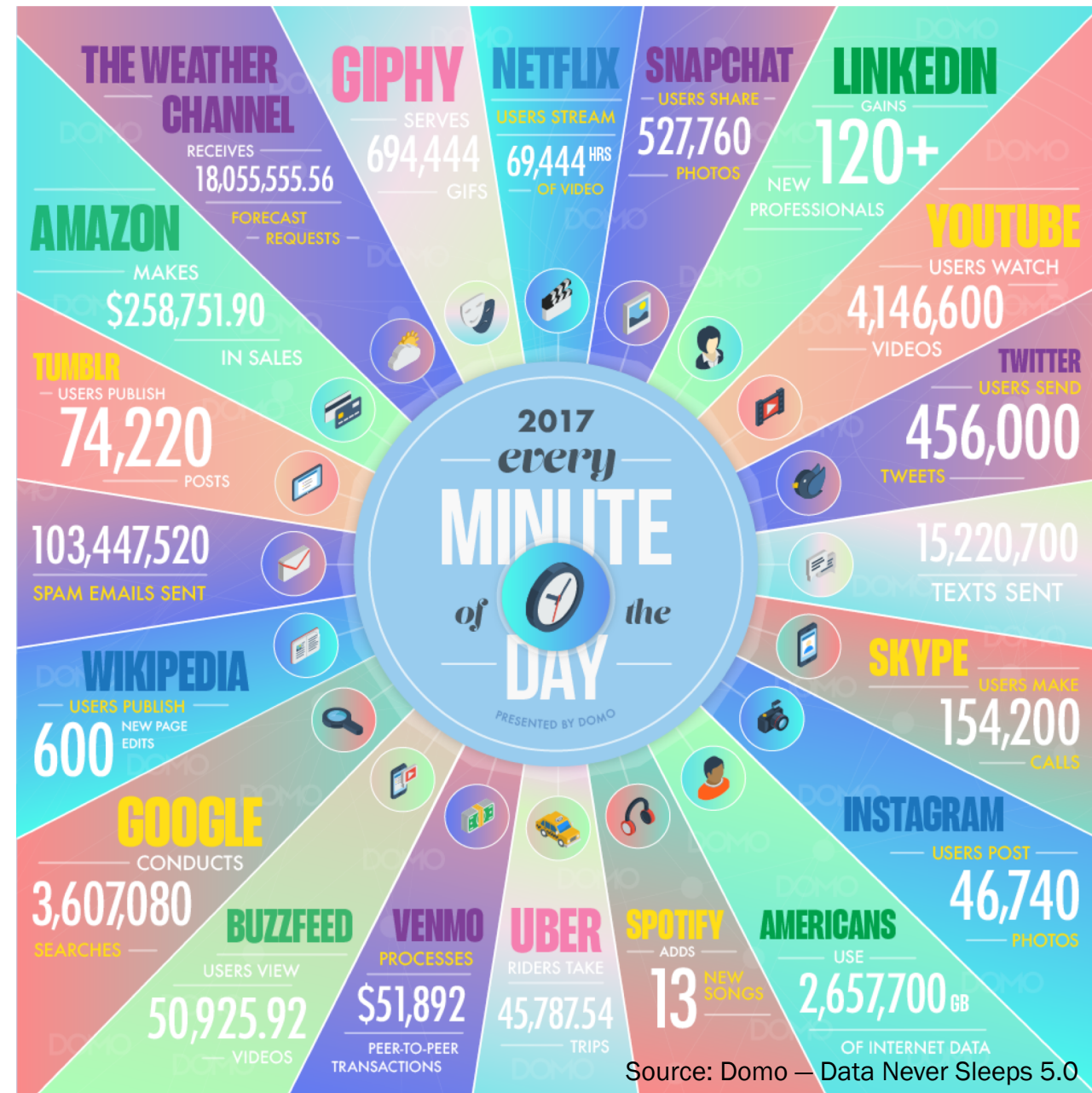\*Computer Science and Engineering, University of Michigan
[†]Department of Computer Science, University of Virginia

UNIVERSITY *of* VIRGINIA

COMPUTER SCIENCE & ENGINEERING
UNIVERSITY OF MICHIGAN

# Processing Growing Quantities of Data

- 2.5 quintillion bytes of data/day
- Analysis /manipulation requires deserialization
- Most data use recursively-nested grammars
  - XML
  - JSON
- Poor performance on CPU
  - High branching



Source: Domo — Data Never Sleeps 5.0

```
<course>
  <footnote></footnote>
  <sln>10637</sln>
  <prefix>ACCTG</prefix>
  <crs>230</crs>
  <lab></lab>
  <sect>01</sect>
  <title>INT FIN ACCT</title>
  <credit>3.0</credit>
  <days>TU,TH</days>
  <times>
     <start>7:45</start>
     <end>9</end>
  </times>
  <place>
      <bldg>TODD</bldg>
      <room>230</room>
  </place>
  <instructor>
    B. MCELDOWNEY
  </instructor>
  <limit>0112</limit>
  <enrolled>0108</enrolled>
</course>
```
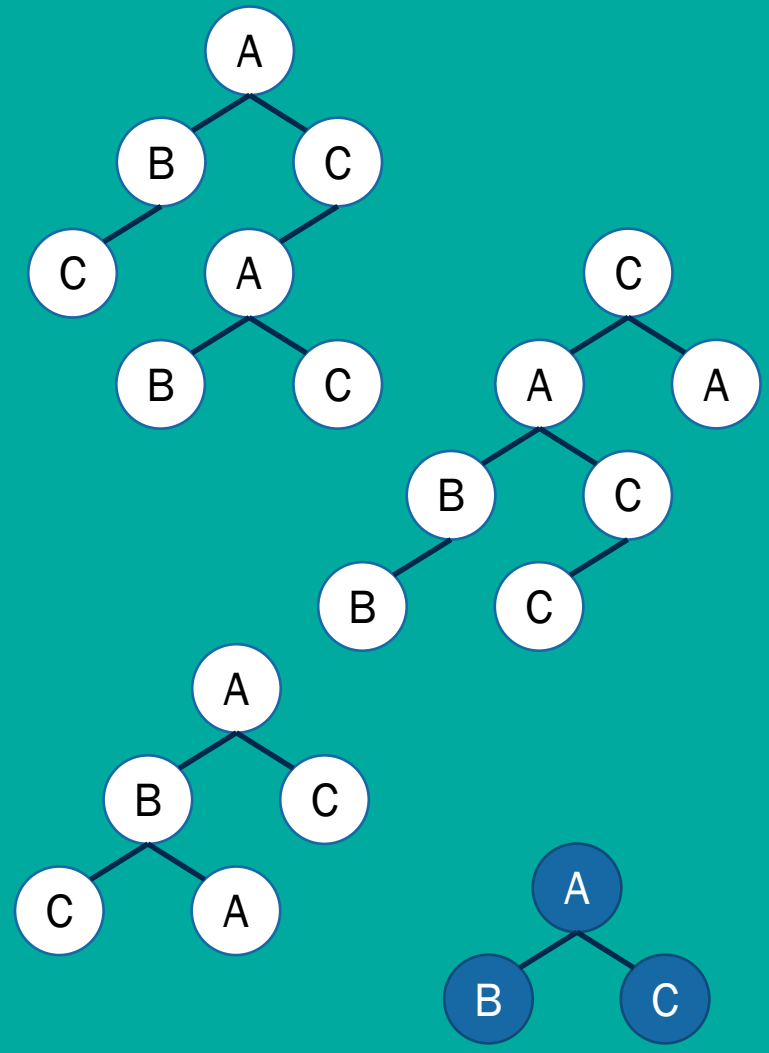
# XML Nesting

```
<course>
  <footnote></footnote>
  <sln>10637</sln>
  <prefix>ACCTG</prefix>
  <crs>230</crs>
  <lab></lab>
  <sect>01</sect>
  <title>INT FIN ACCT</title>
  <credit>3.0</credit>
  <days>TU,TH</days>
  <times>
      <start>7:45</start>
      <end>9</end>
  </times>
  <place>
      <bldg>TODD</bldg>
      <room>230</room>
  </place>
  <instructor>
    B. MCELDOWNEY
  </instructor>
  <limit>0112</limit>
  <enrolled>0108</enrolled>
</course>
```
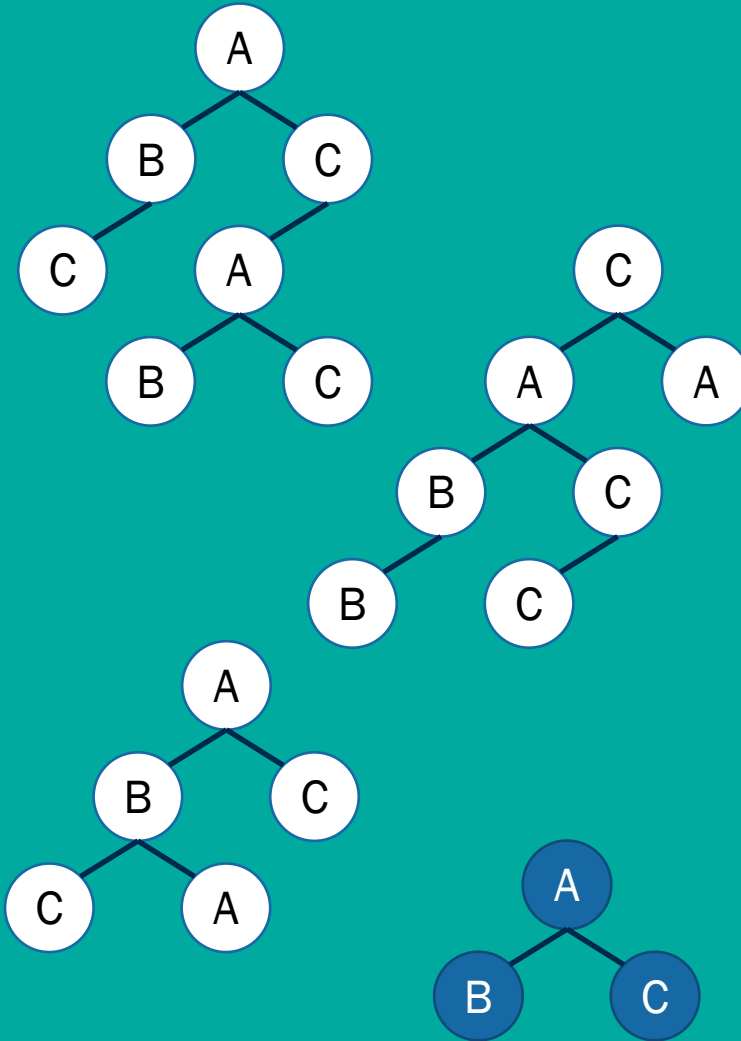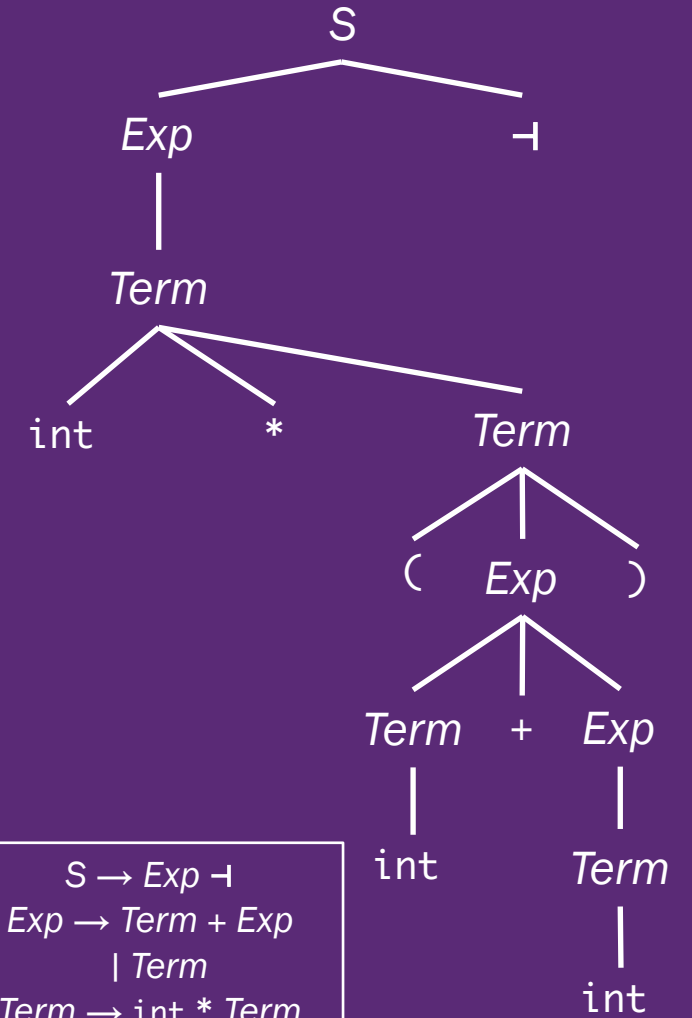
**XML Nesting**

**Subtree Mining**

**XML Nesting**

```
<course>
  <footnote></footnote>
  <sln>10637</sln>
  <prefix>ACCTG</prefix>
  <crs>230</crs>
  <lab></lab>
  <sect>01</sect>
  <title>INT FIN ACCT</title>
  <credit>3.0</credit>
  <days>TU,TH</days>
  <times>
      <start>7:45</start>
      <end>9</end>
  </times>
  <place>
      <bldg>TODD</bldg>
      <room>230</room>
  </place>
  <instructor>
    B. MCELDOWNEY
  </instructor>
  <limit>0112</limit>
  <enrolled>0108</enrolled>
</course>
```
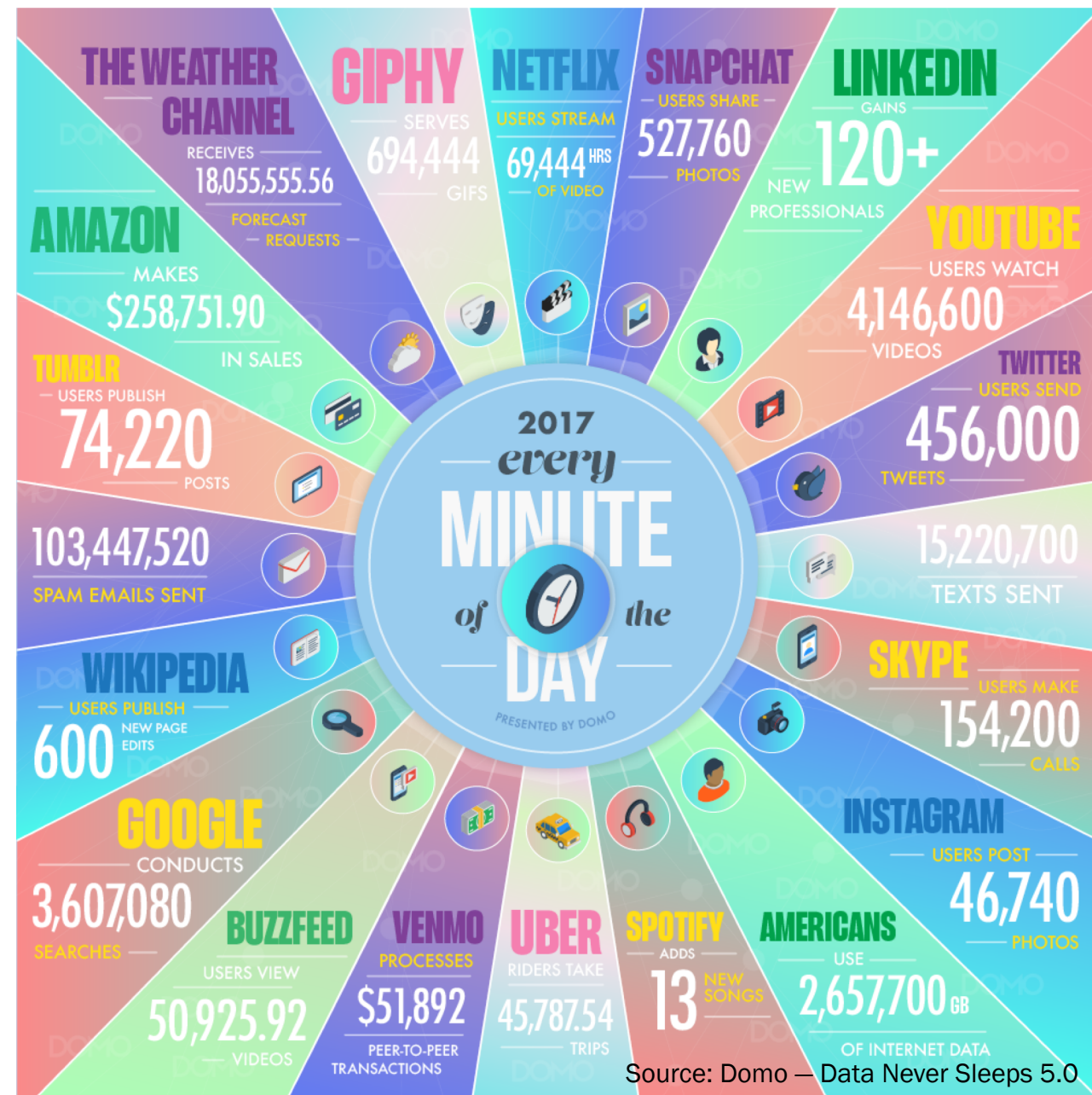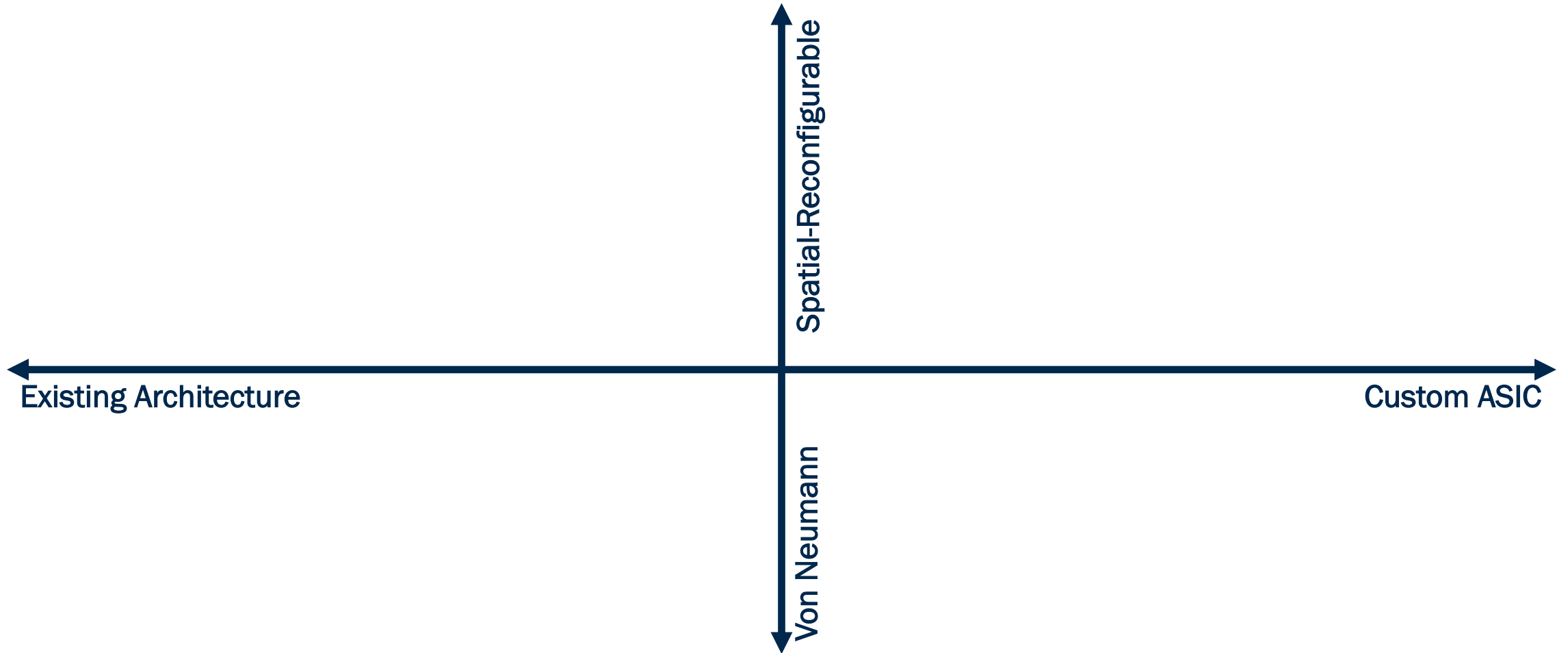
**Subtree Mining**

**Parsing**

$S \rightarrow Exp \dashv$
$Exp \rightarrow Term + Exp$
$\qquad | Term$
$Term \rightarrow \text{int} * Term$
$\qquad | ( Exp )$
$\qquad | \text{int}$
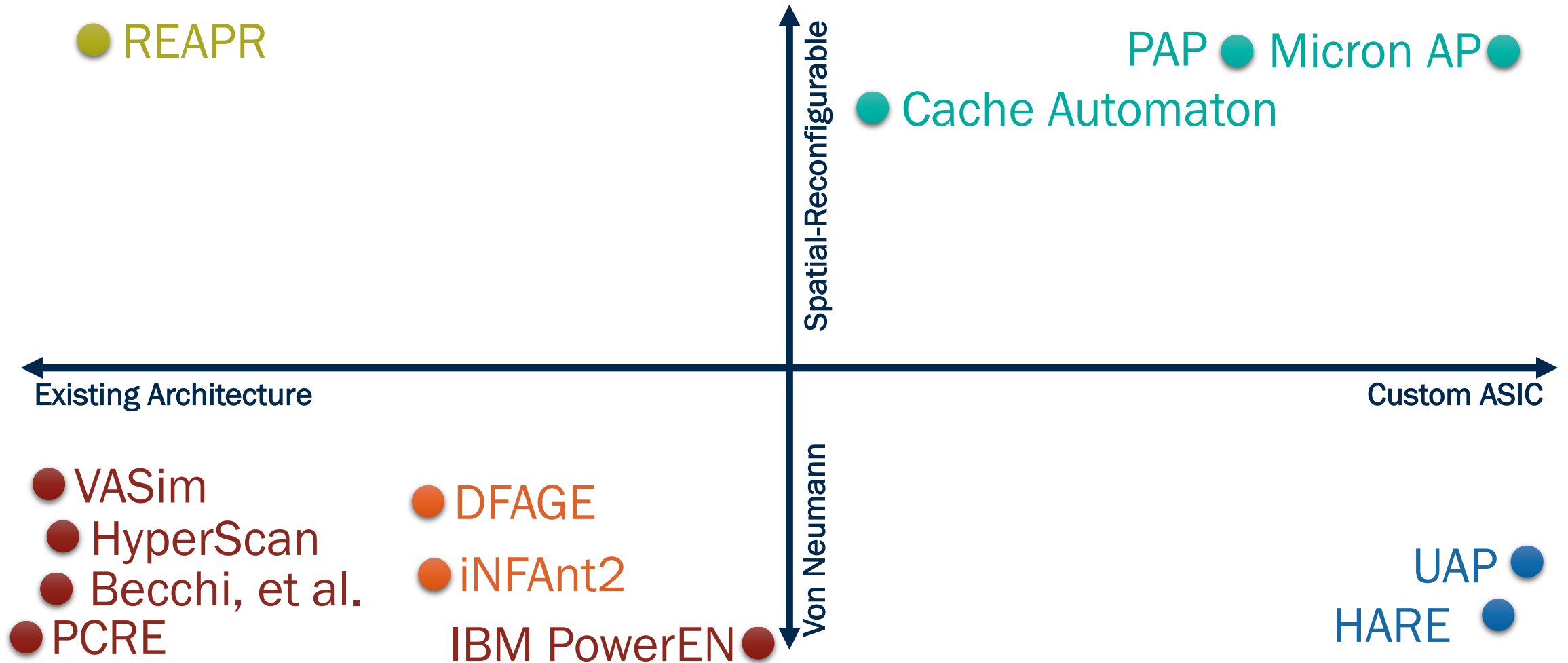
# Processing Growing Quantities of Data

- Automata/RegEx help tame analysis of big data sets
  - Frequent Itemset/Pattern Mining
  - NLP Part-of-Speech Tagging
  - Data Deduplication
  - Ensemble-Based Classification
  - Particle Physics Analyses
- Growing number of architectural solutions



Source: Domo — Data Never Sleeps 5.0

# Automata/RegEx Processing Platforms



Spatial-Reconfigurable

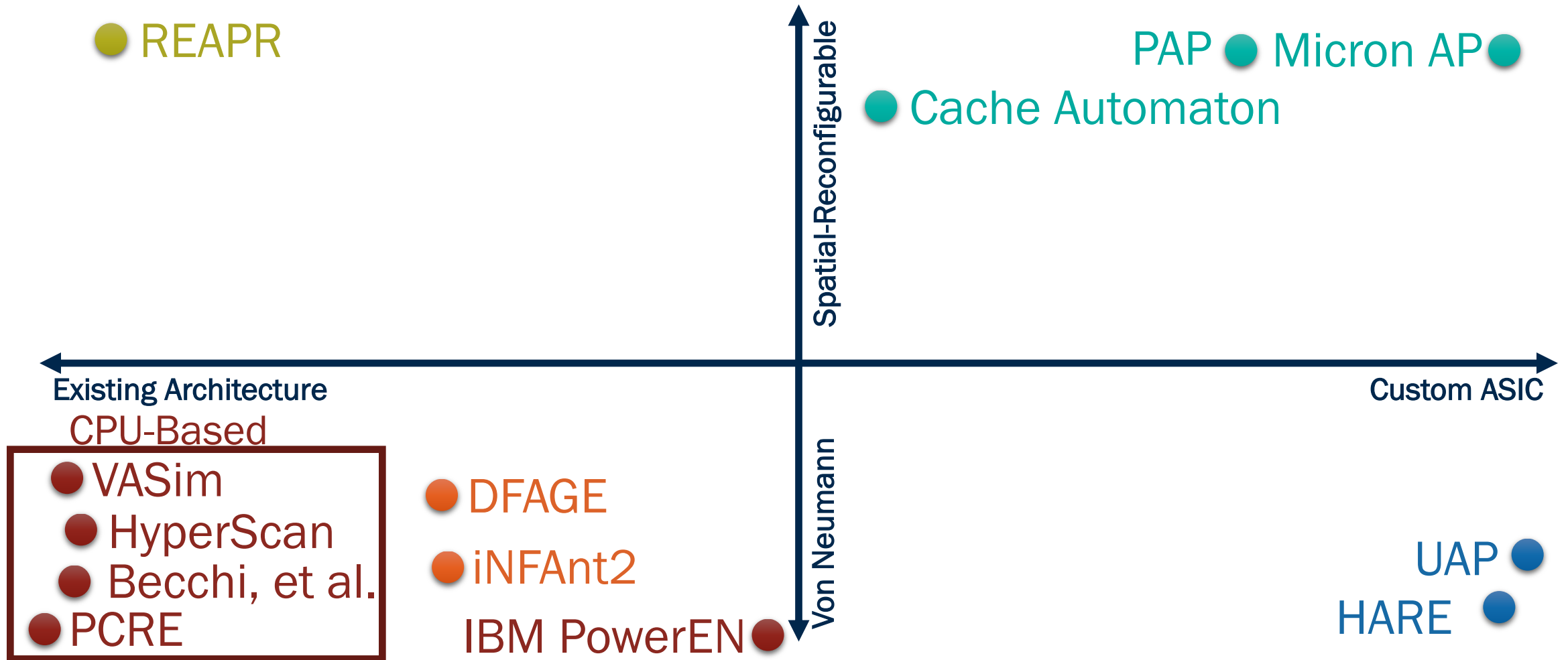Existing Architecture ← → Custom ASIC

Von Neumann
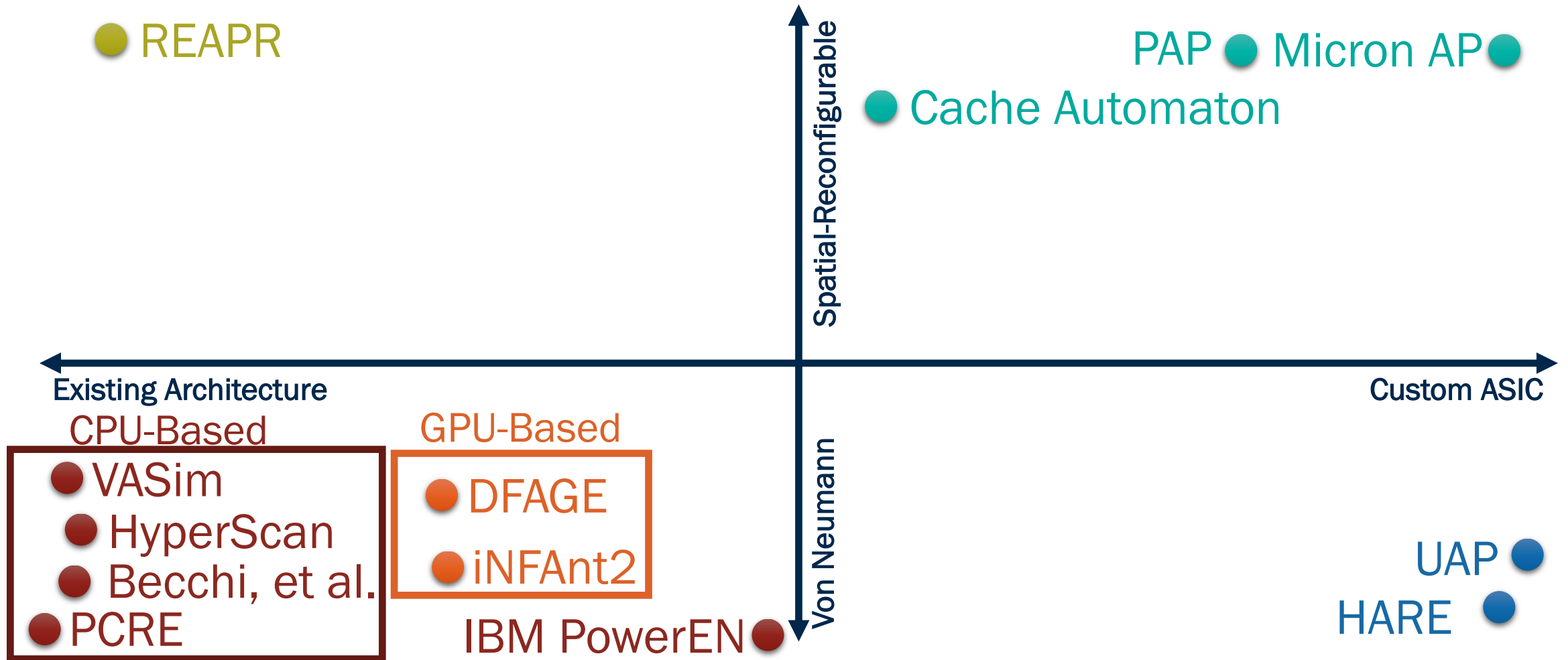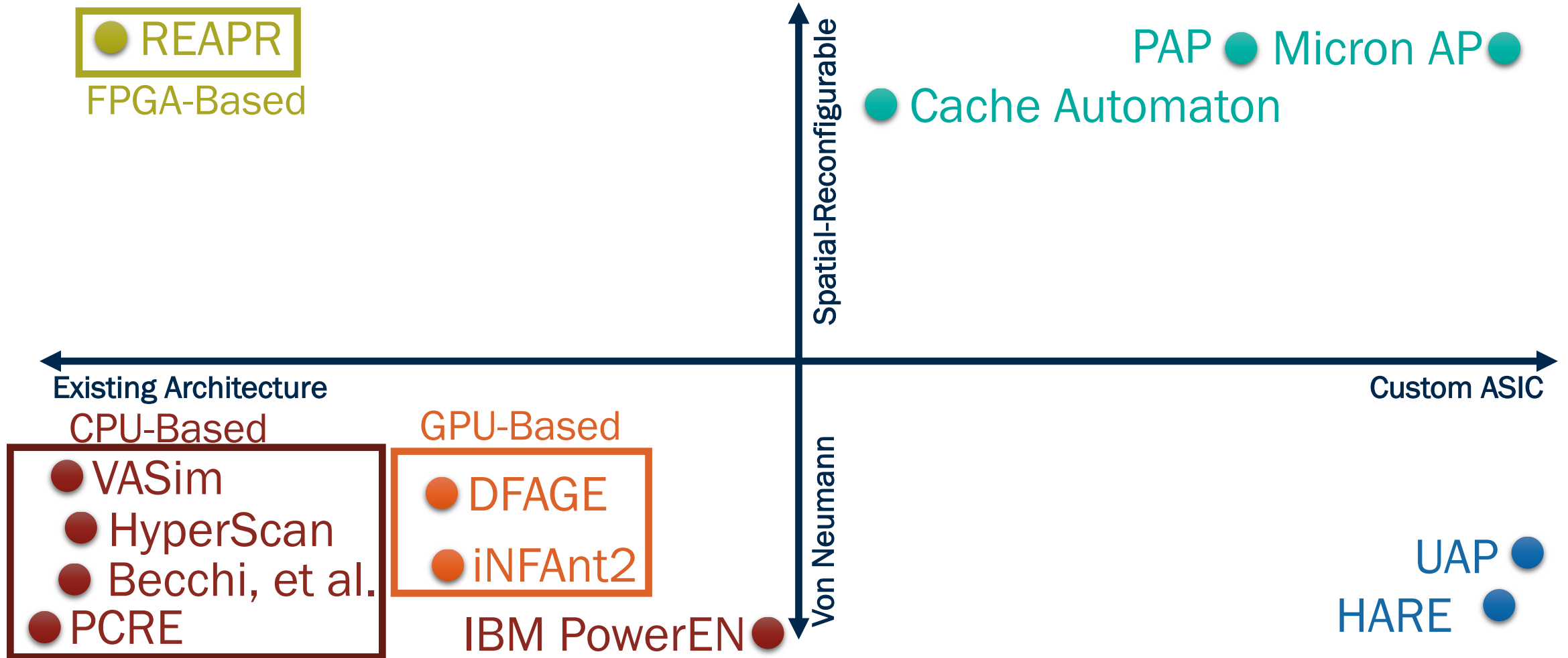
# Automata/RegEx Processing Platforms

# Automata/RegEx Processing Platforms

# Automata/RegEx Processing Platforms



Spatial-Reconfigurable

REAPR
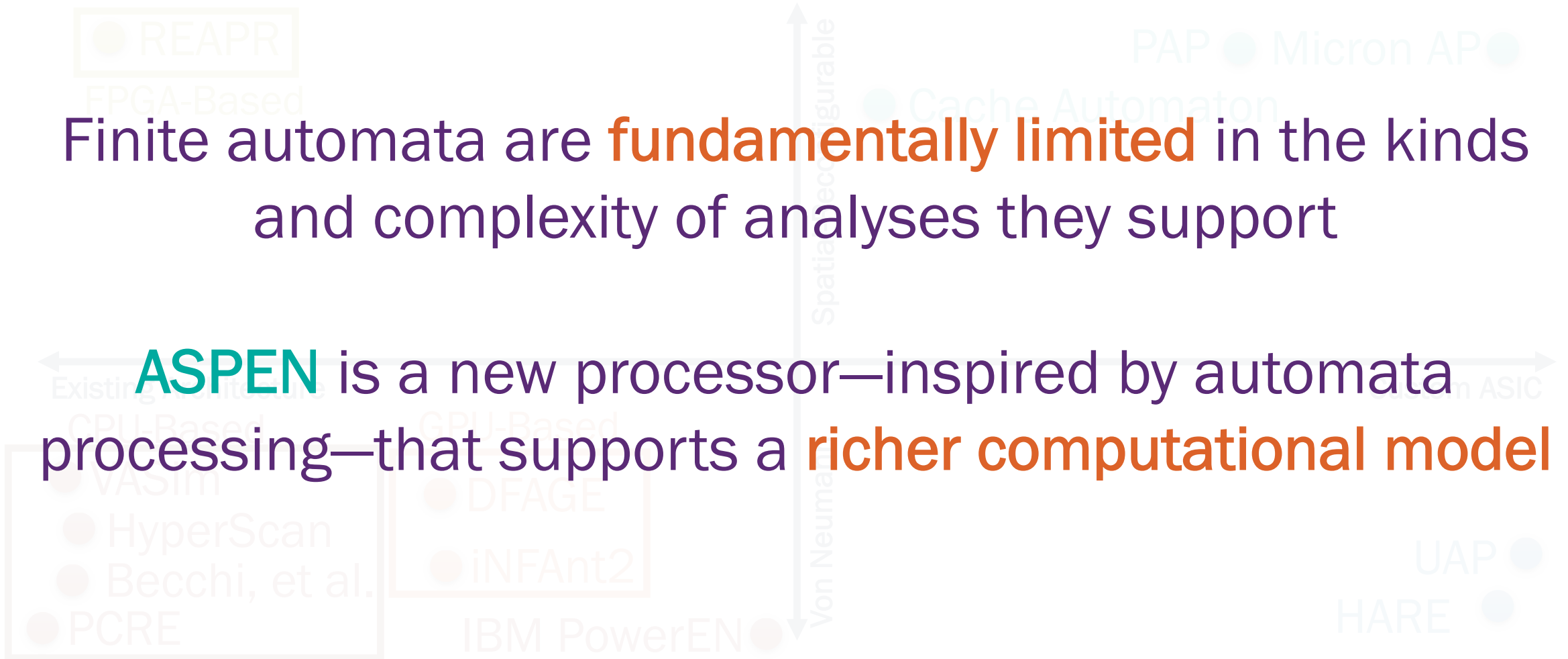
PAP • Micron AP •

• Cache Automaton

Existing Architecture

Custom ASIC

CPU-Based

GPU-Based

Von Neumann

• VASim
• HyperScan
• Becchi, et al.
• PCRE

• DFAGE
• iNFAnt2

IBM PowerEN •

UAP •

HARE •

# Automata/RegEx Processing Platforms

# Automata/RegEx Processing Platforms

Finite automata are fundamentally limited in the kinds and complexity of analyses they support

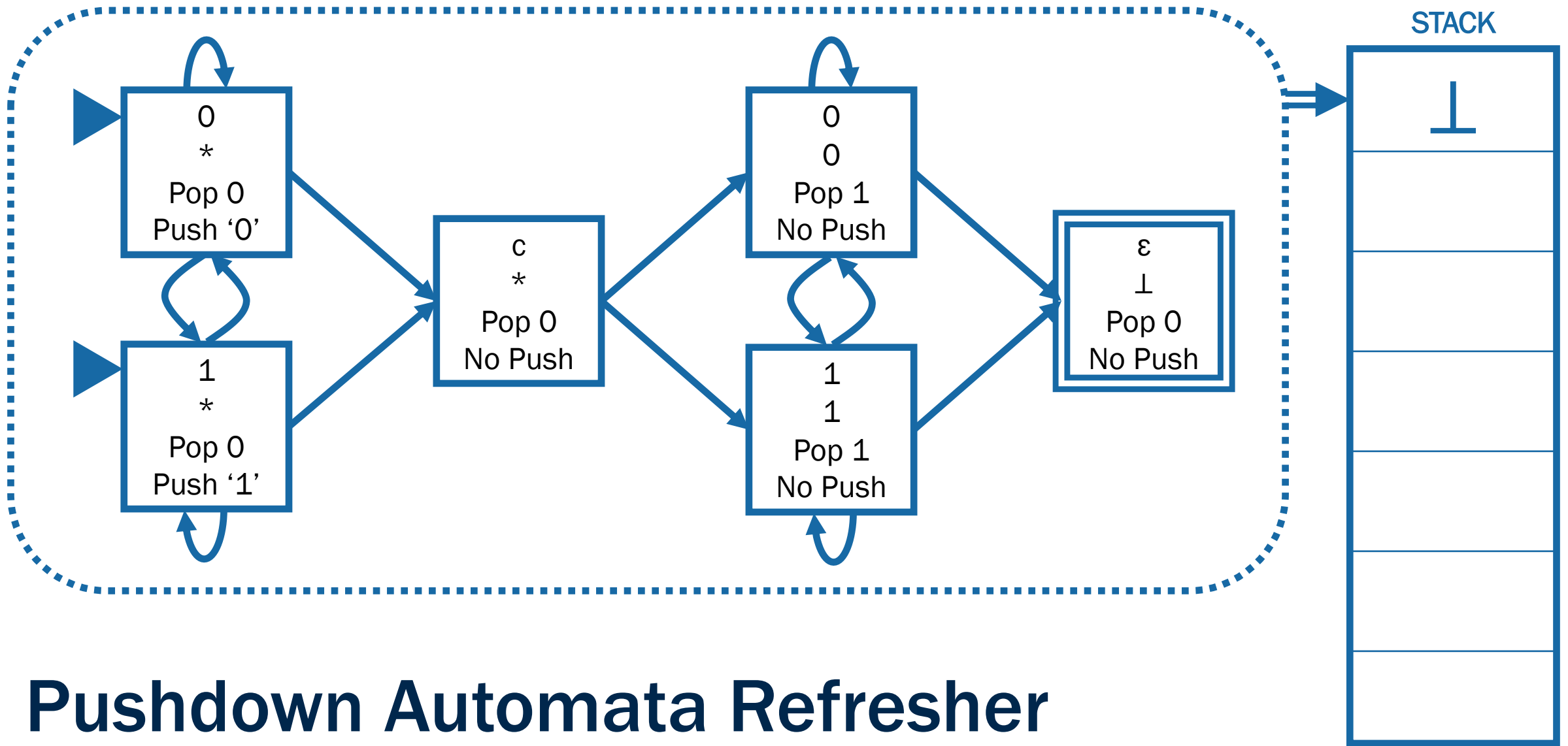ASPEN is a new processor—inspired by automata processing—that supports a richer computational model

# Automata/RegEx Processing Platforms

Finite automata are **fundamentally limited** in the kinds and complexity of analyses they support

ASPEN is a new processor—inspired by automata processing—that supports a **richer computational model**

# ASPEN Supports Richer Analyses

- **A**ccelerated in-**S**RAM **P**ushdown **EN**gine

- Scalable processing engine that uses LLC slices to accelerate Pushdown Automata computation

- Custom five-stage datapath using SRAM lookups can process up to one byte per cycle

- Optimizing compiler supports existing grammars, packs states efficiently, and reduces the number processing stalls

- Provides additional cache when not in use

# Overview of this Talk

- Pushdown Automata Refresher

- Architectural Design of ASPEN
  - Why LLC?
  - Datapath innovation

- Optimizations

- Evaluation
  - XML Parsing
  - Subtree Mining

# Pushdown Automata Refresher
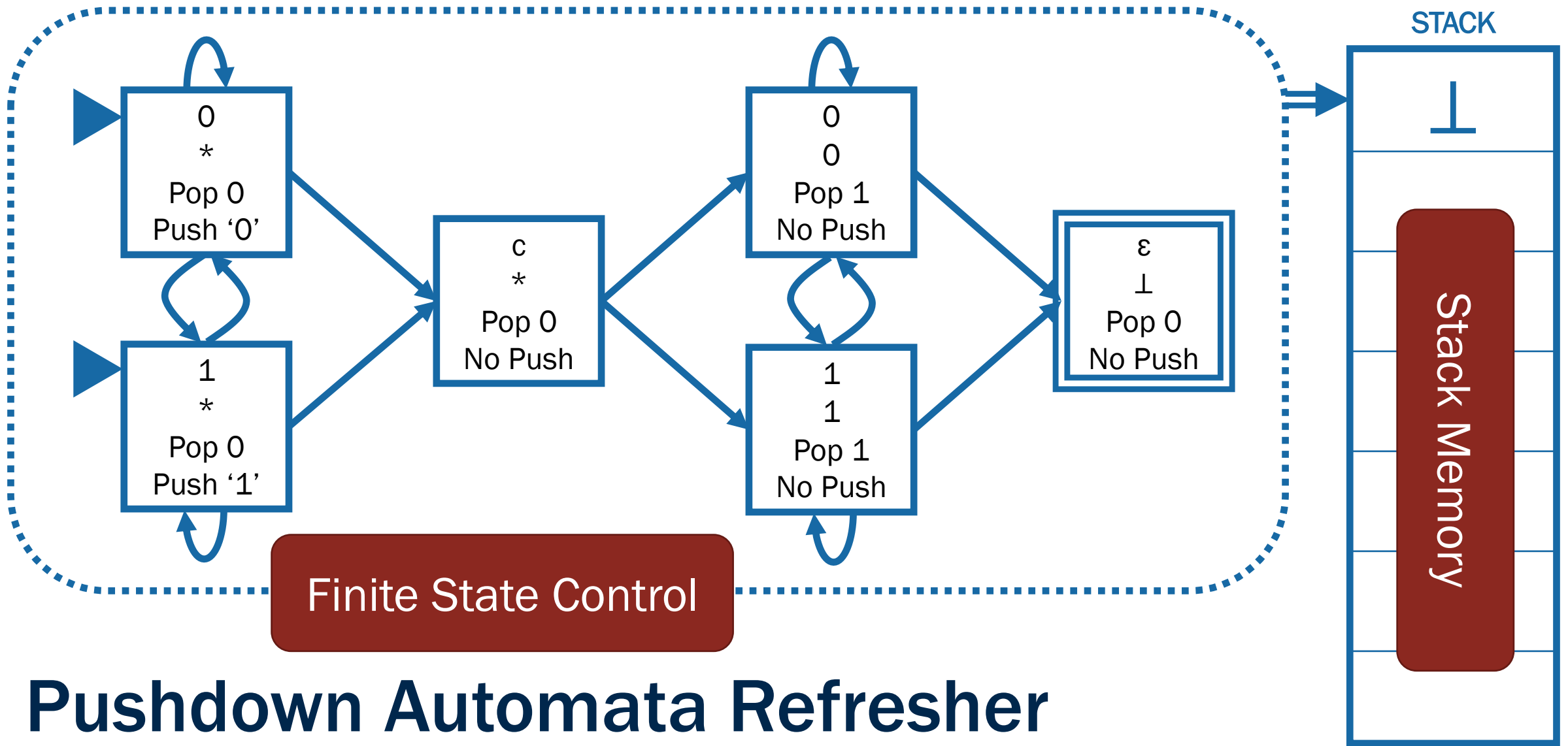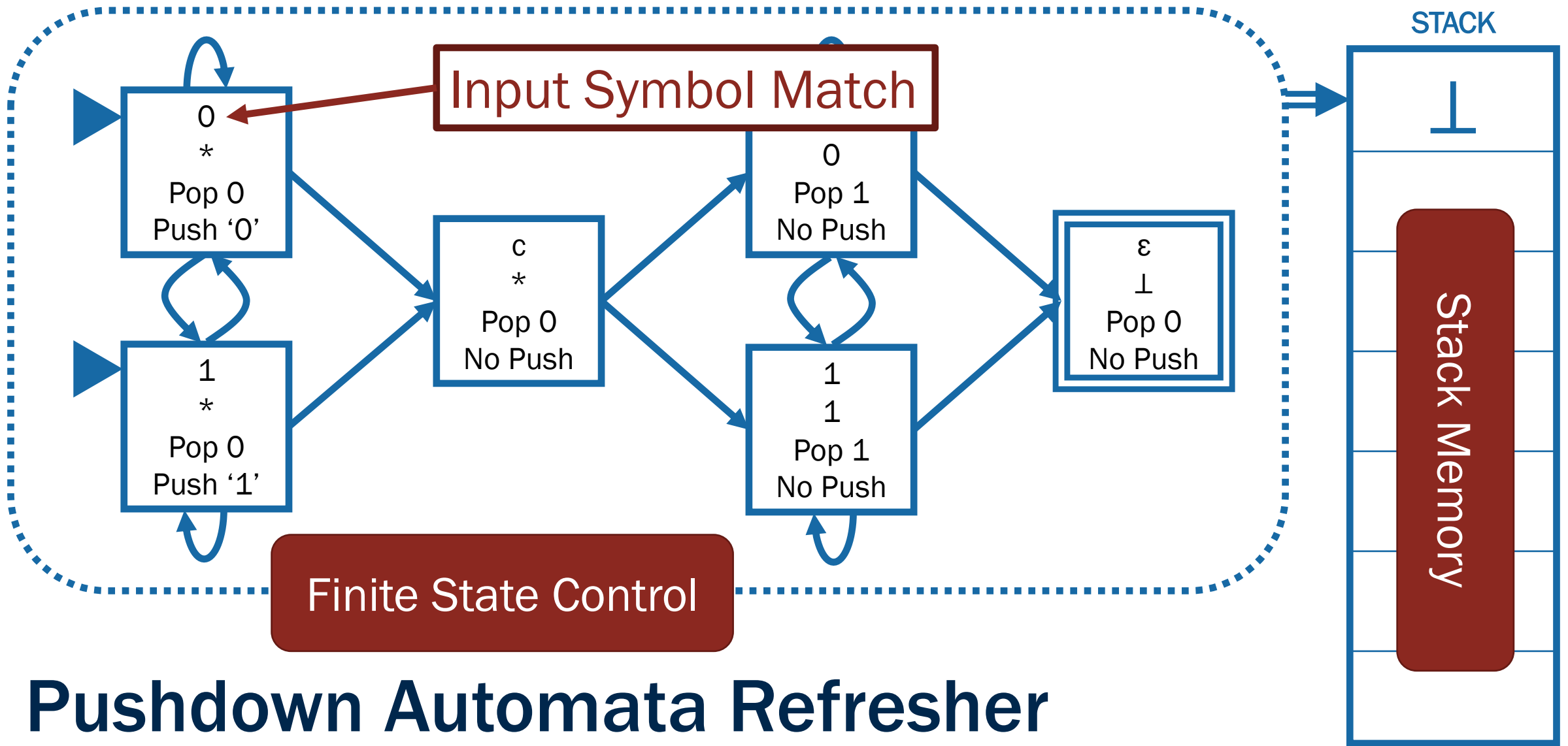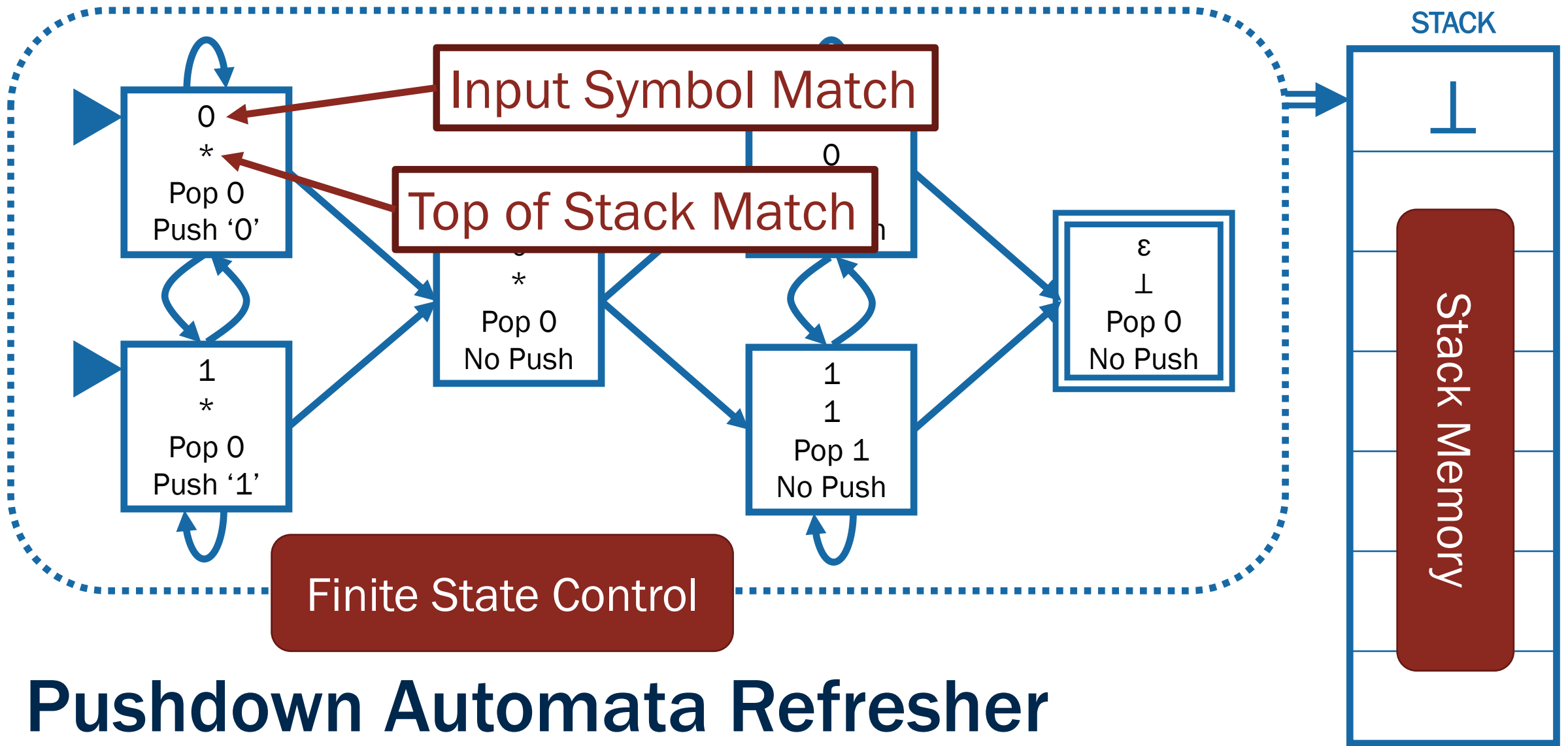
STACK

16

# Pushdown Automata Refresher

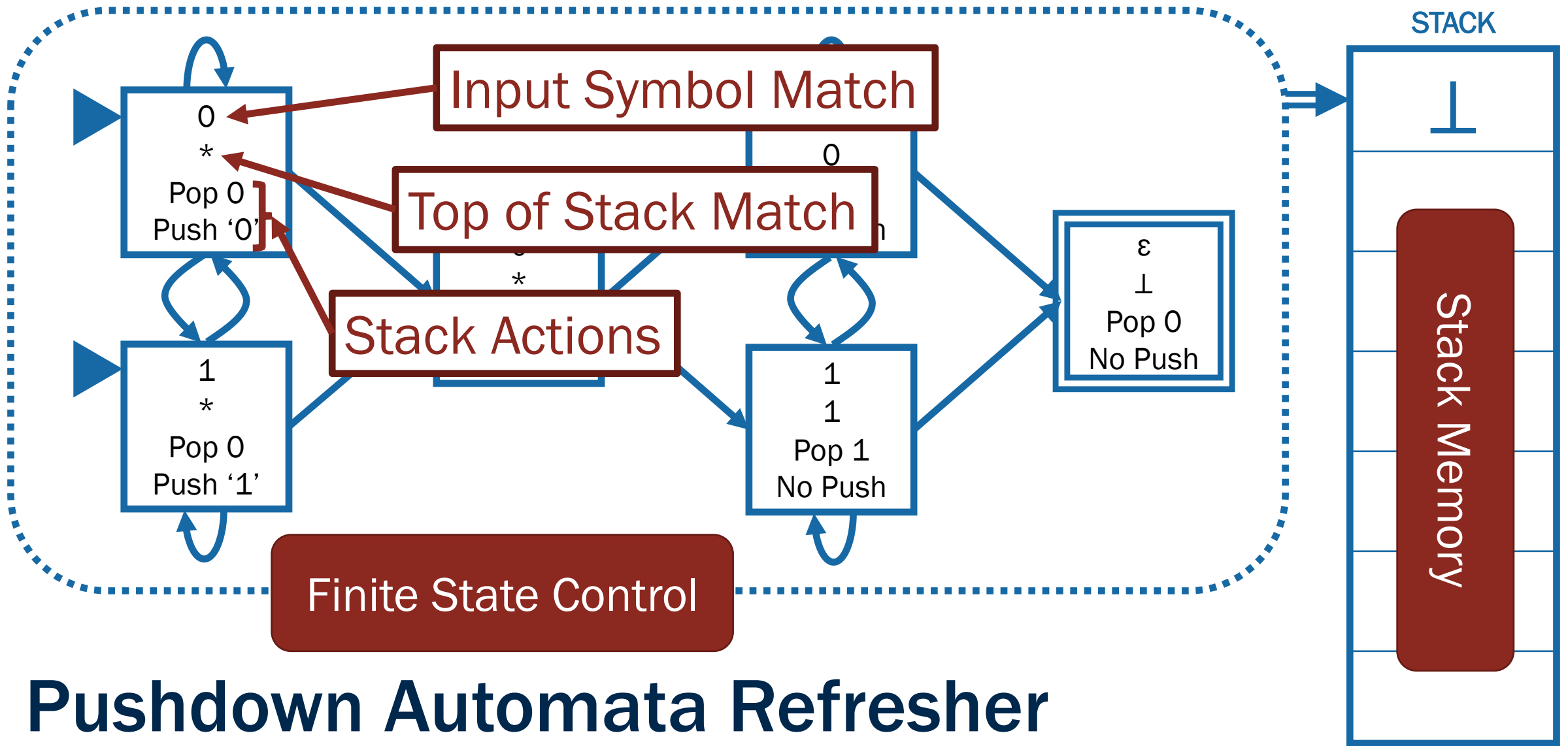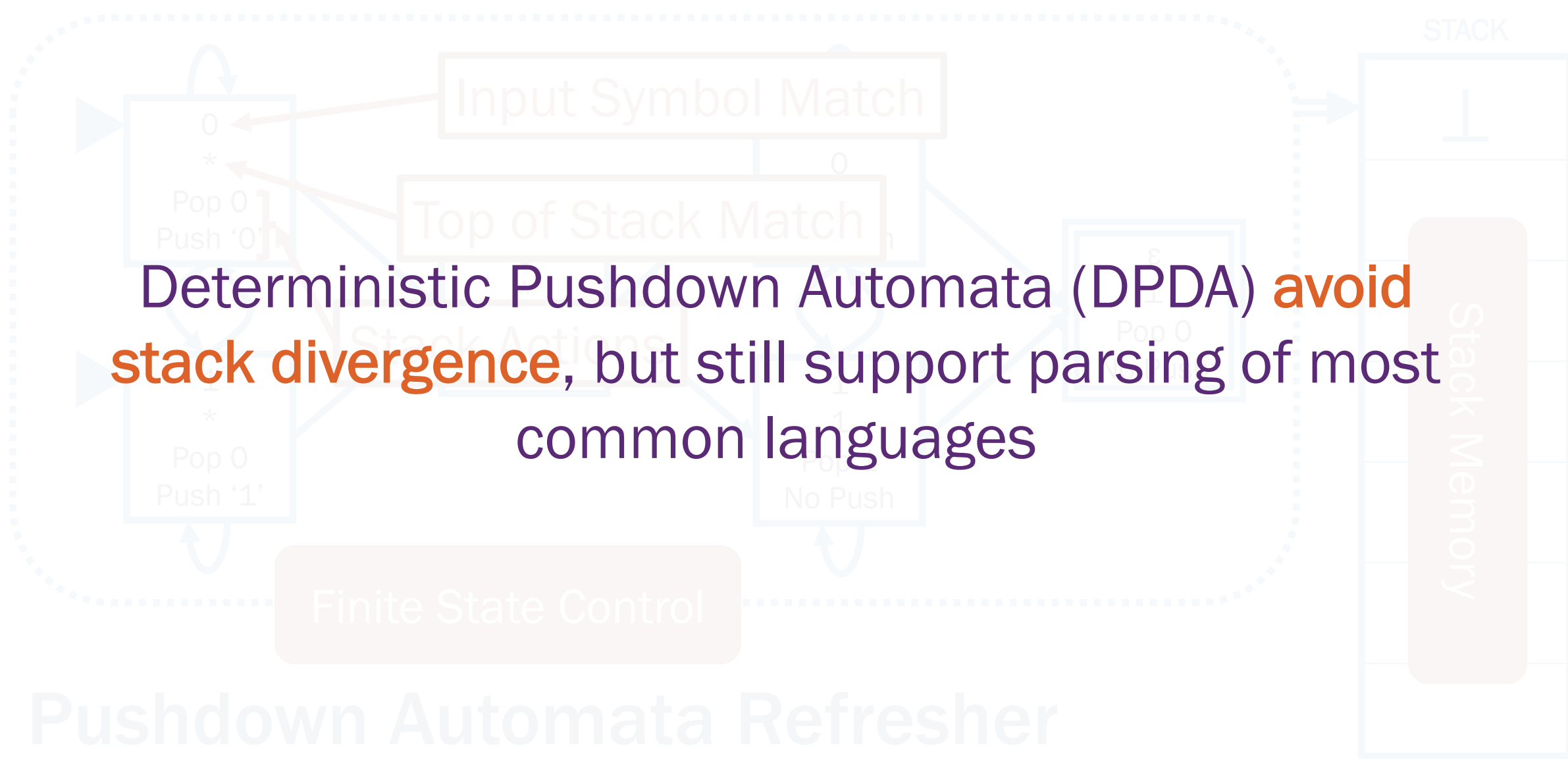Pushdown Automata Refresher

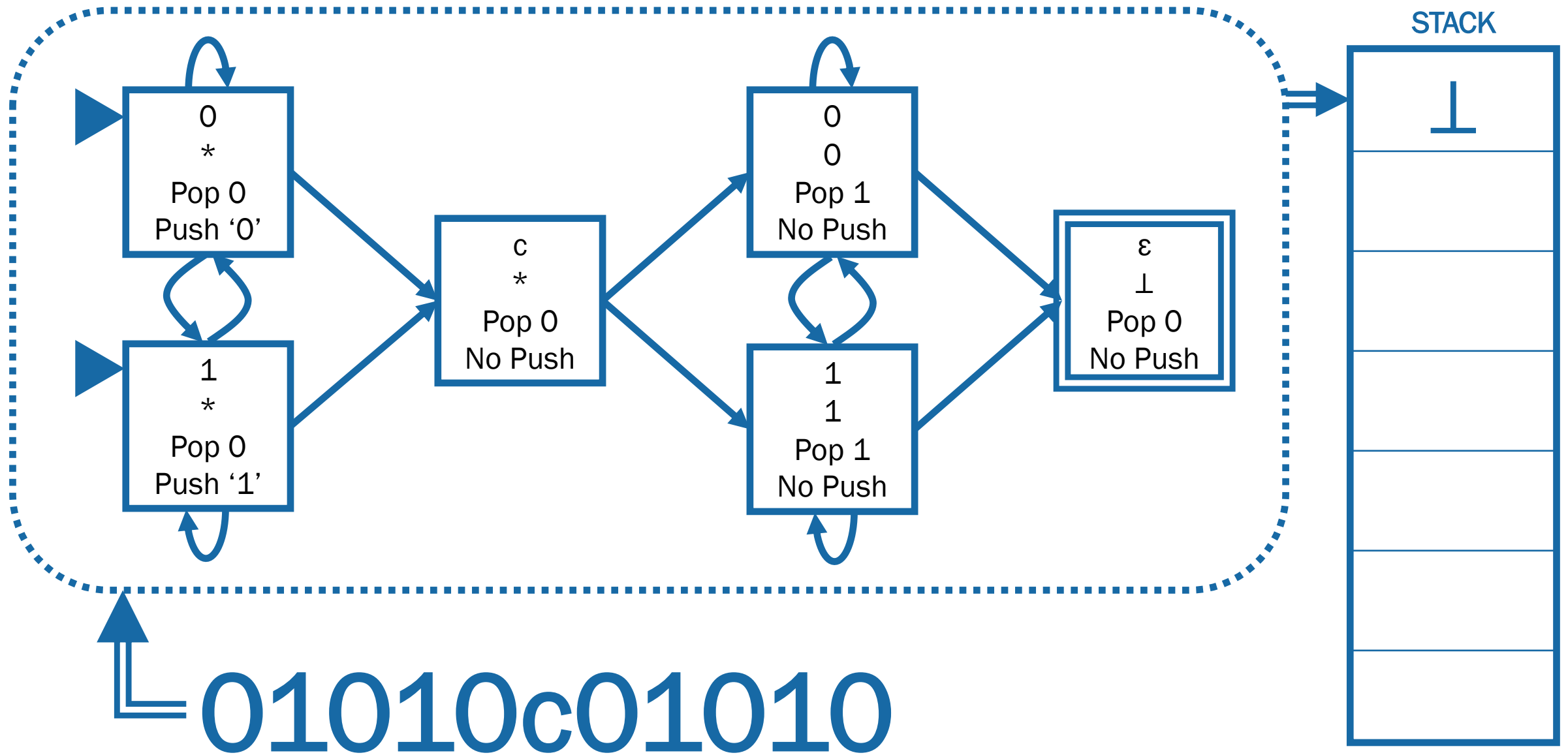Pushdown Automata Refresher

# Pushdown Automata Refresher

**Pushdown Automata Refresher**

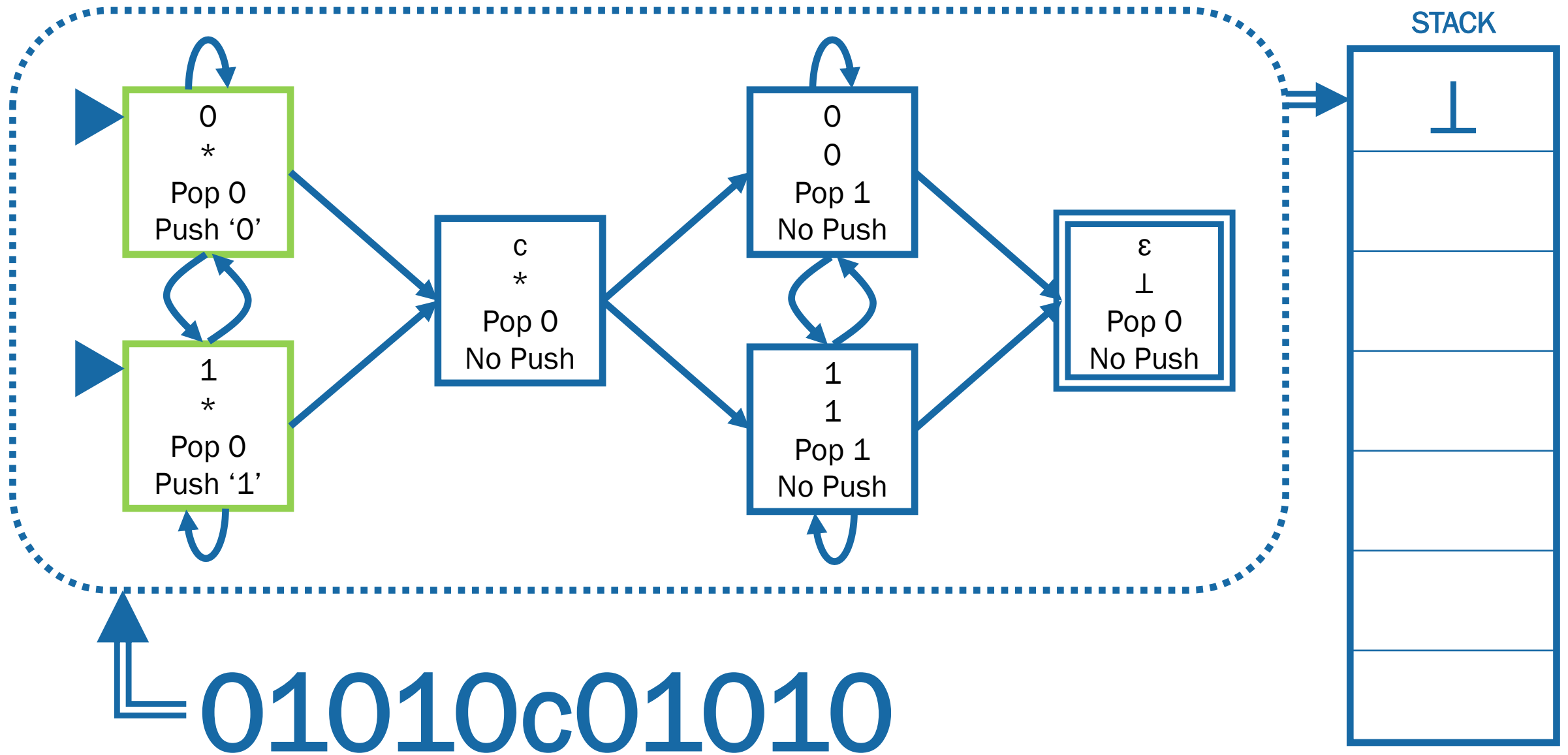Deterministic Pushdown Automata (DPDA) avoid stack divergence, but still support parsing of most common languages

# Recognizing Palindromes with a Middle Character



STACK

01010c01010

23

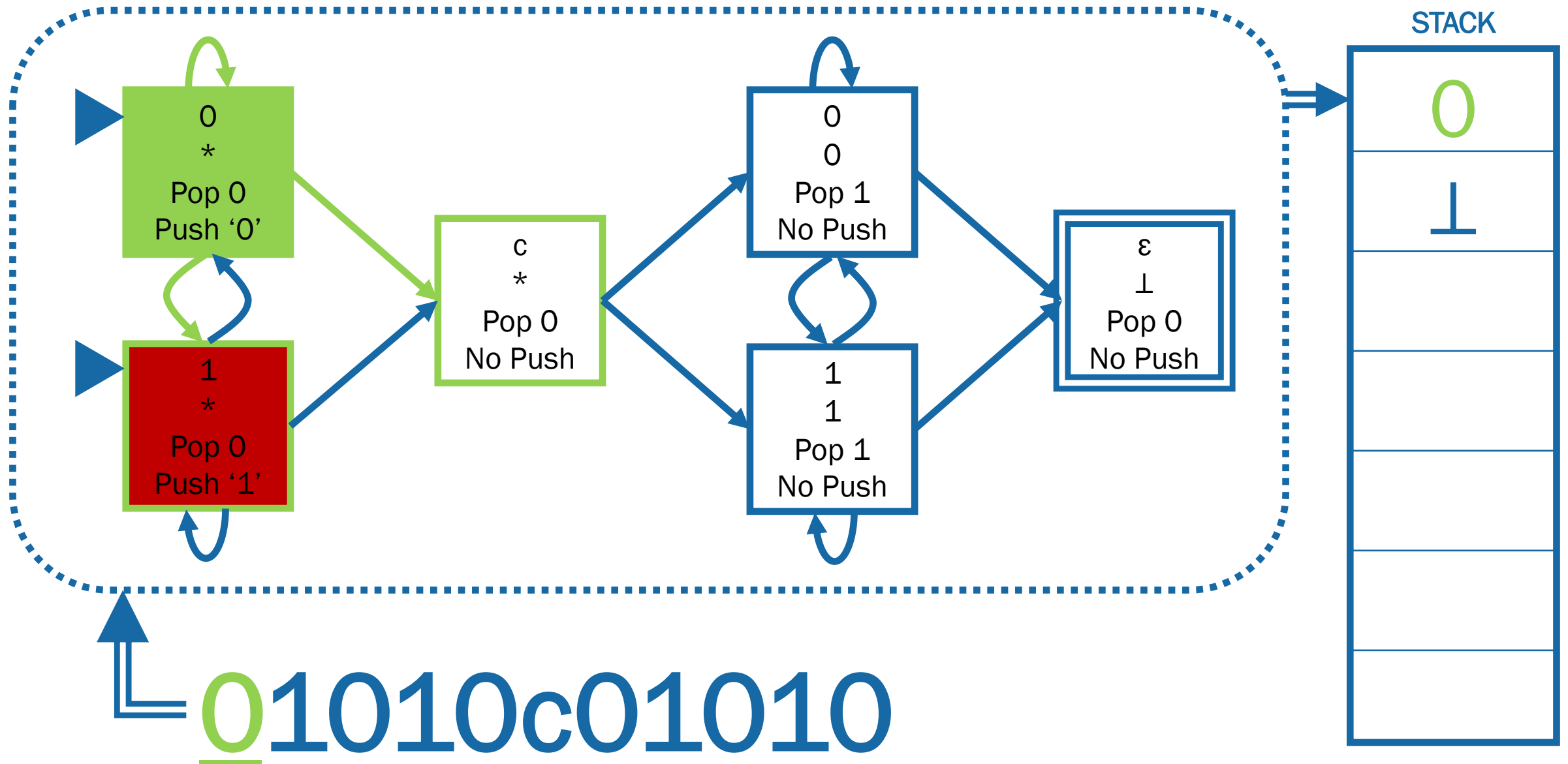# Recognizing Palindromes with a Middle Character



STACK

01010c01010

# Recognizing Palindromes with a Middle Character

# Recognizing Palindromes with a Middle Character

# Recognizing Palindromes with a Middle Character

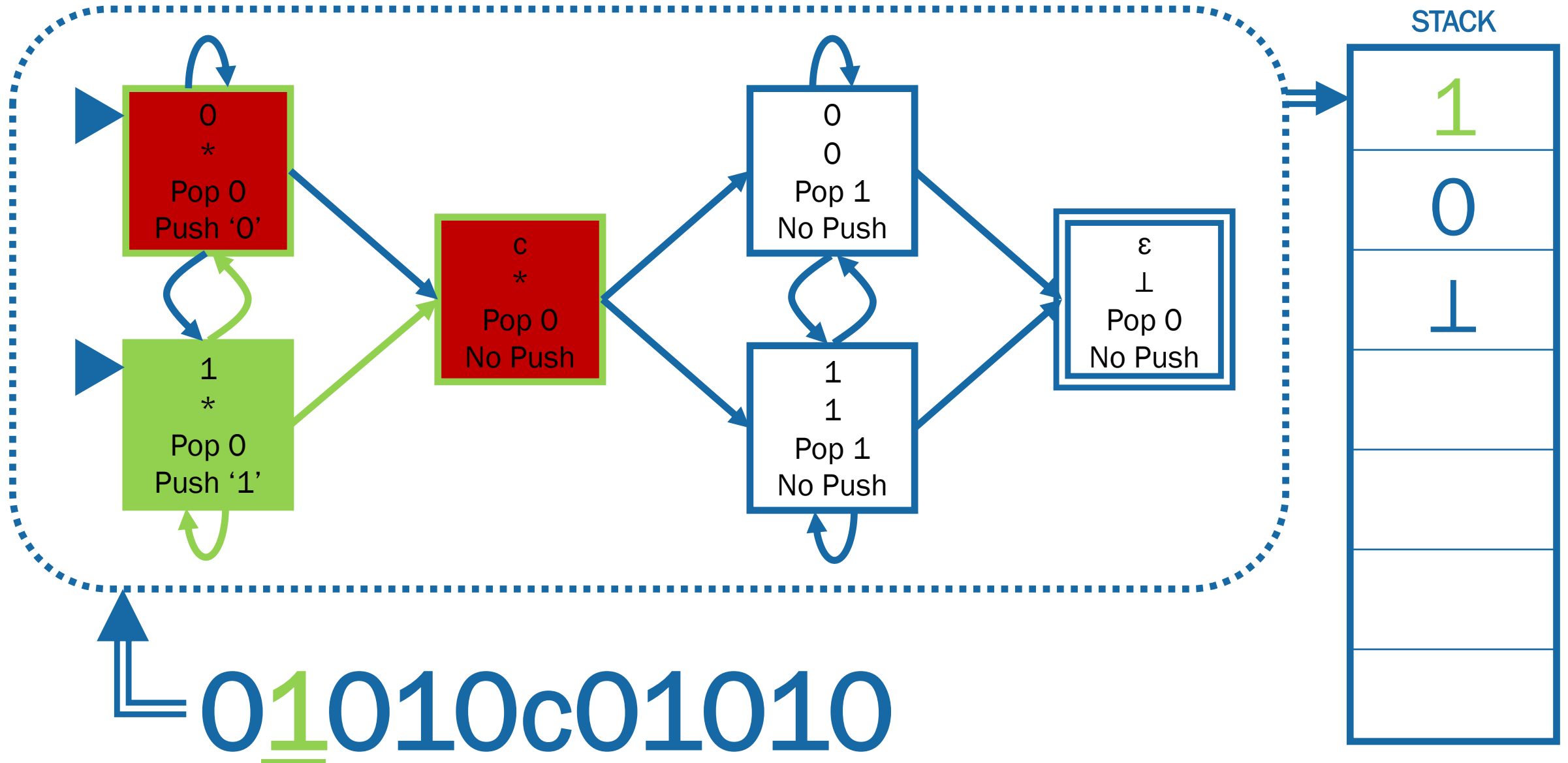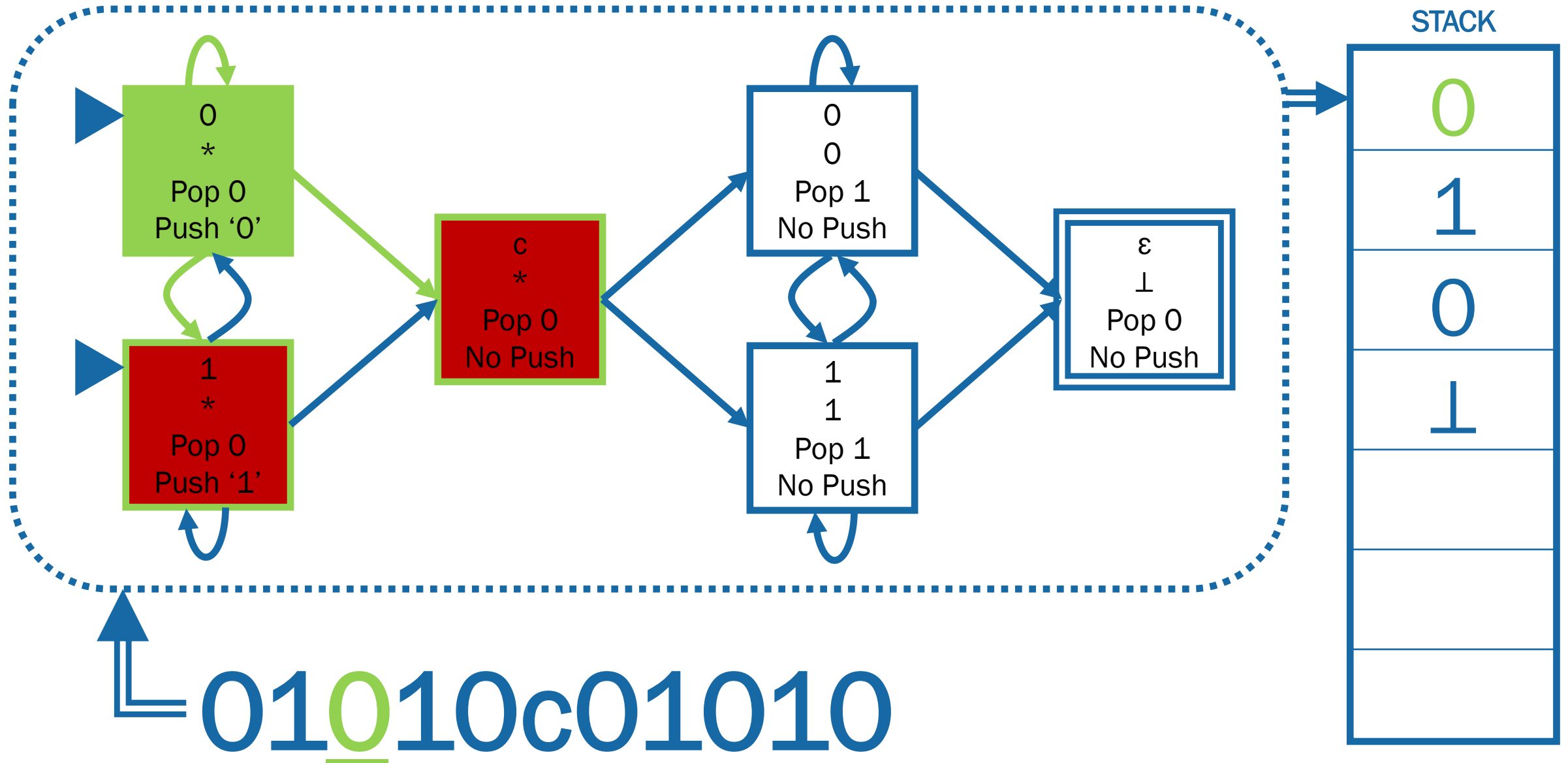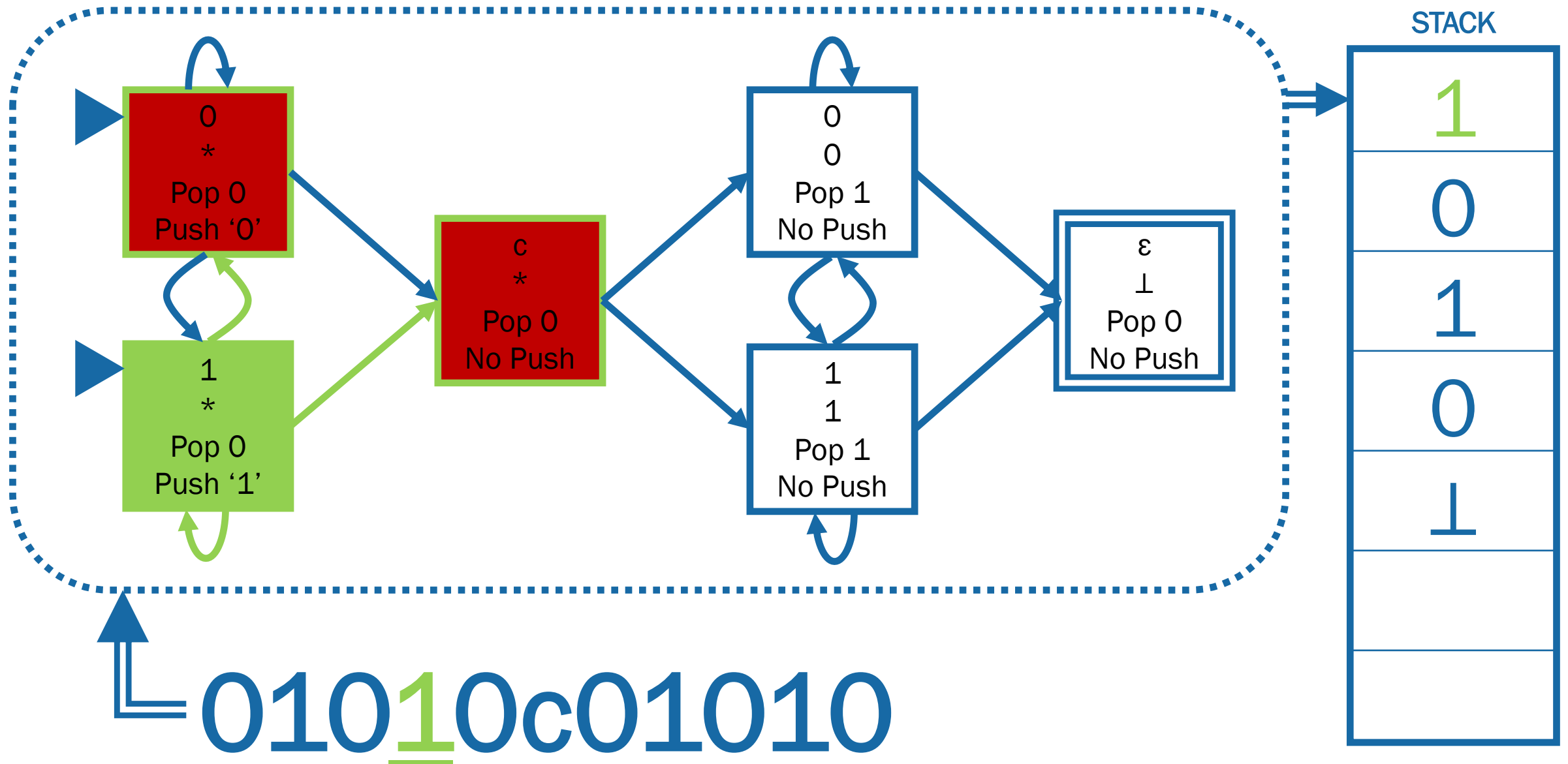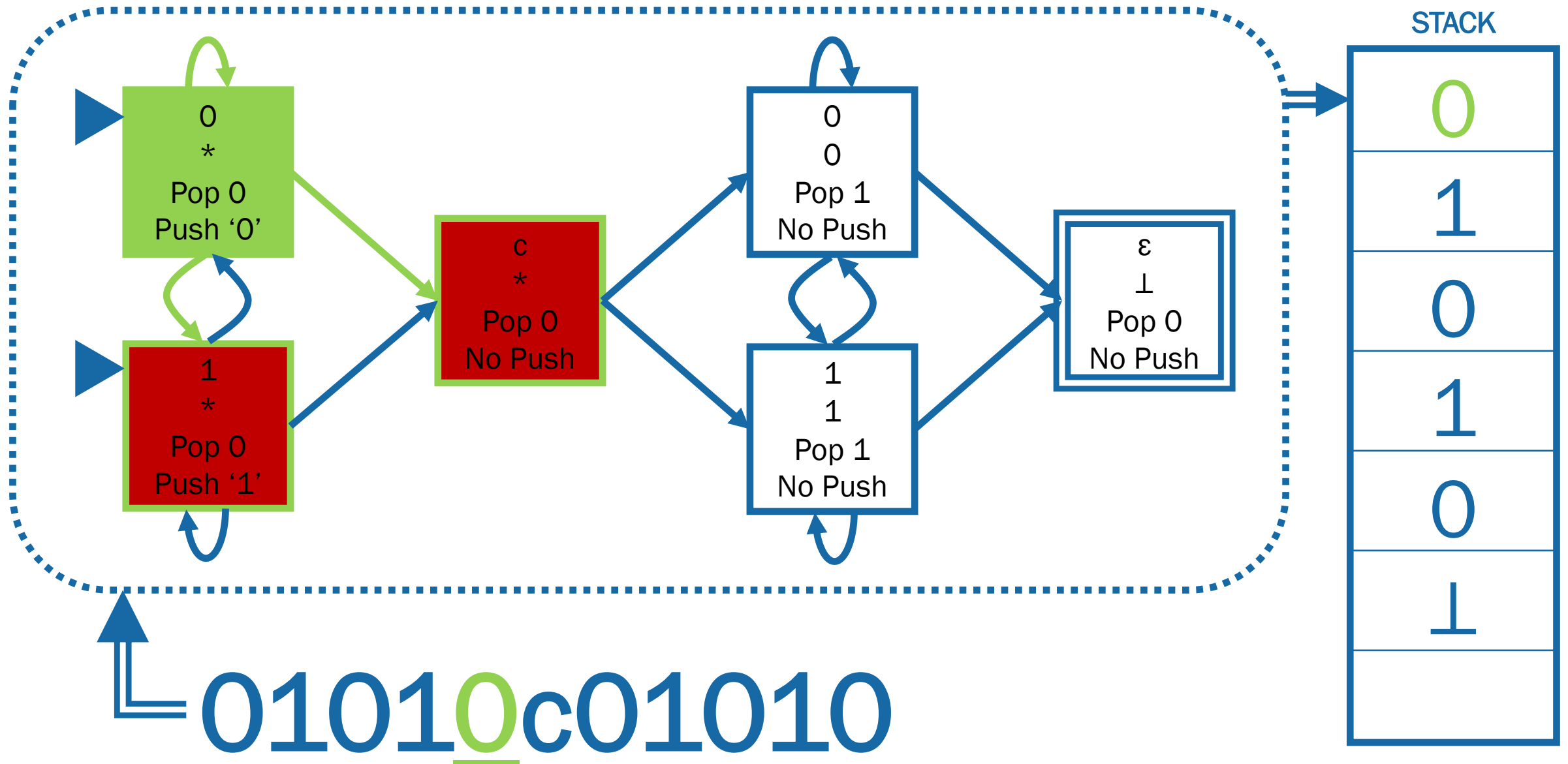# Recognizing Palindromes with a Middle Character

# Recognizing Palindromes with a Middle Character



STACK

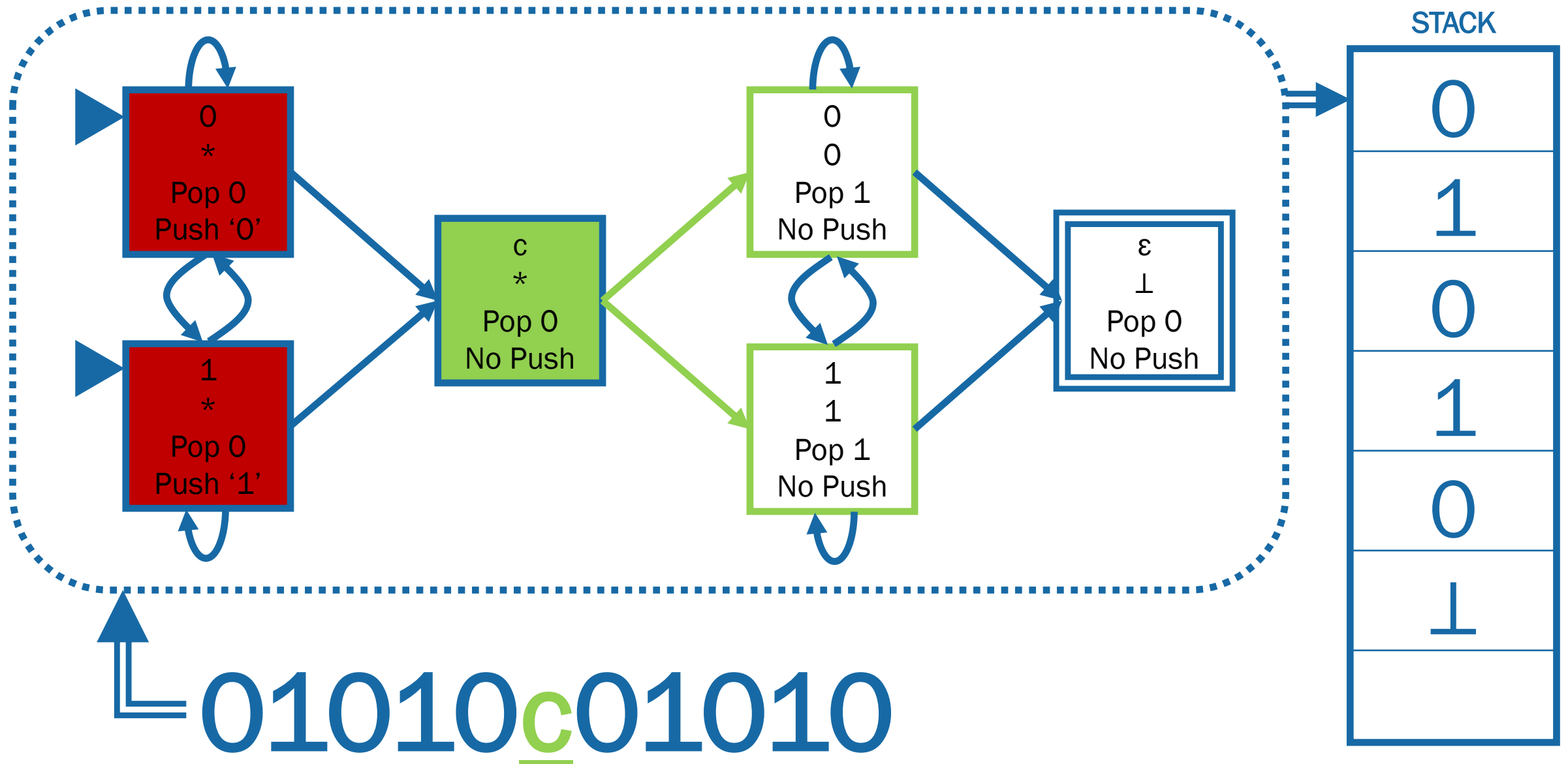0
1
0
1
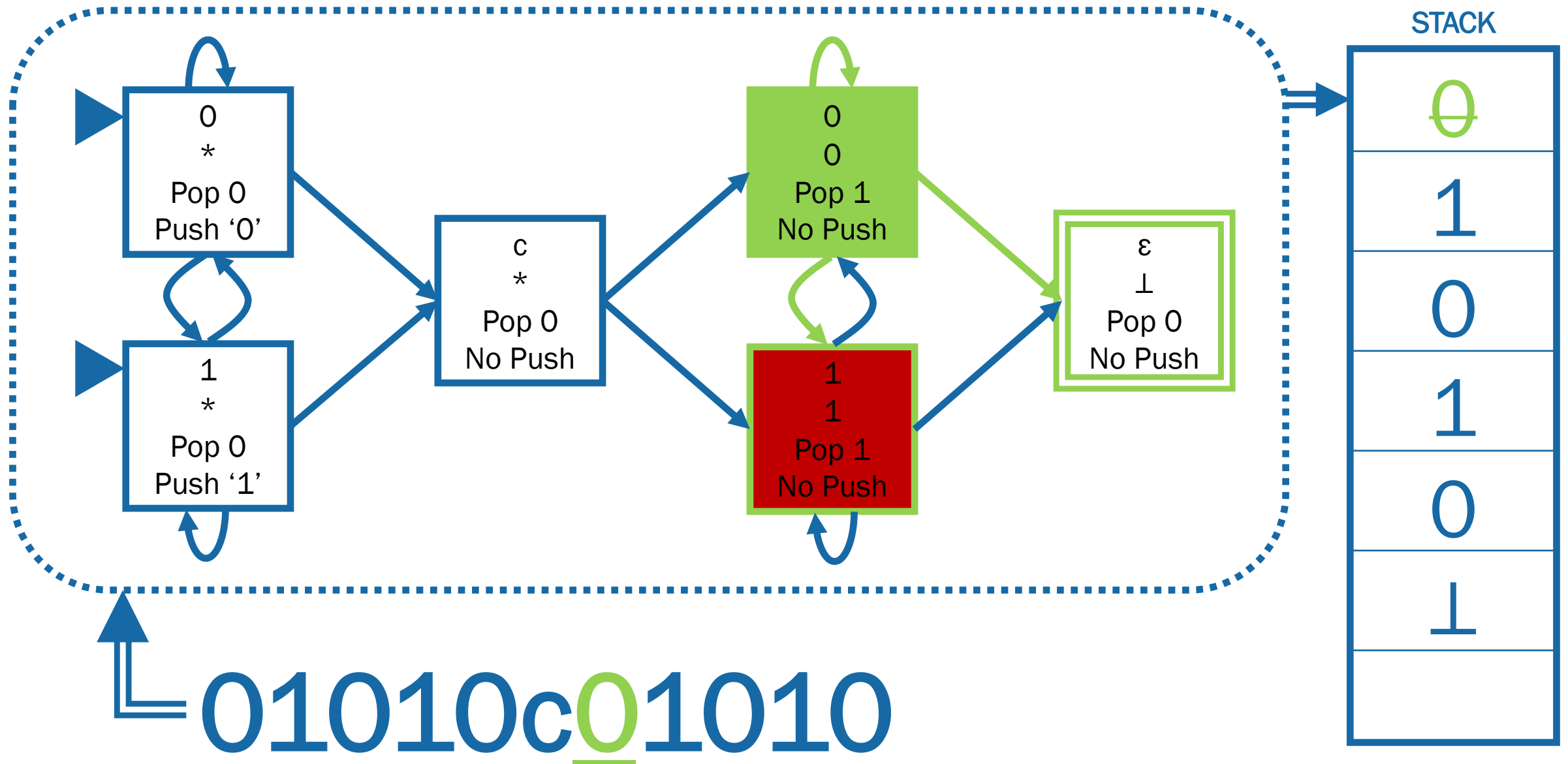0
⊥

01010c01010

29

# Recognizing Palindromes with a Middle Character

# Recognizing Palindromes with a Middle Character



STACK

31

# Recognizing Palindromes with a Middle Character



STACK

01010c01010

32

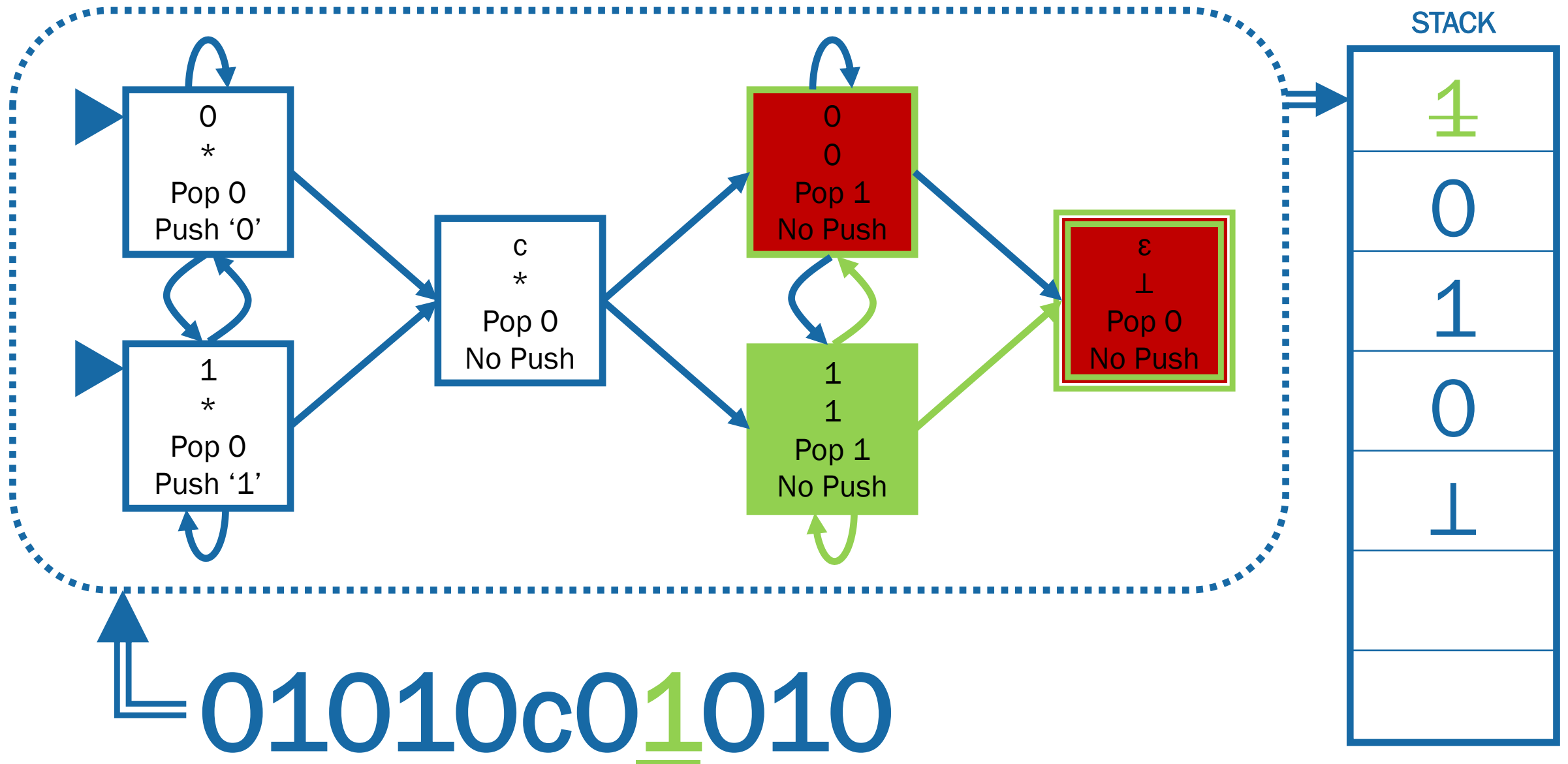# Recognizing Palindromes with a Middle Character



33

# Recognizing Palindromes with a Middle Character
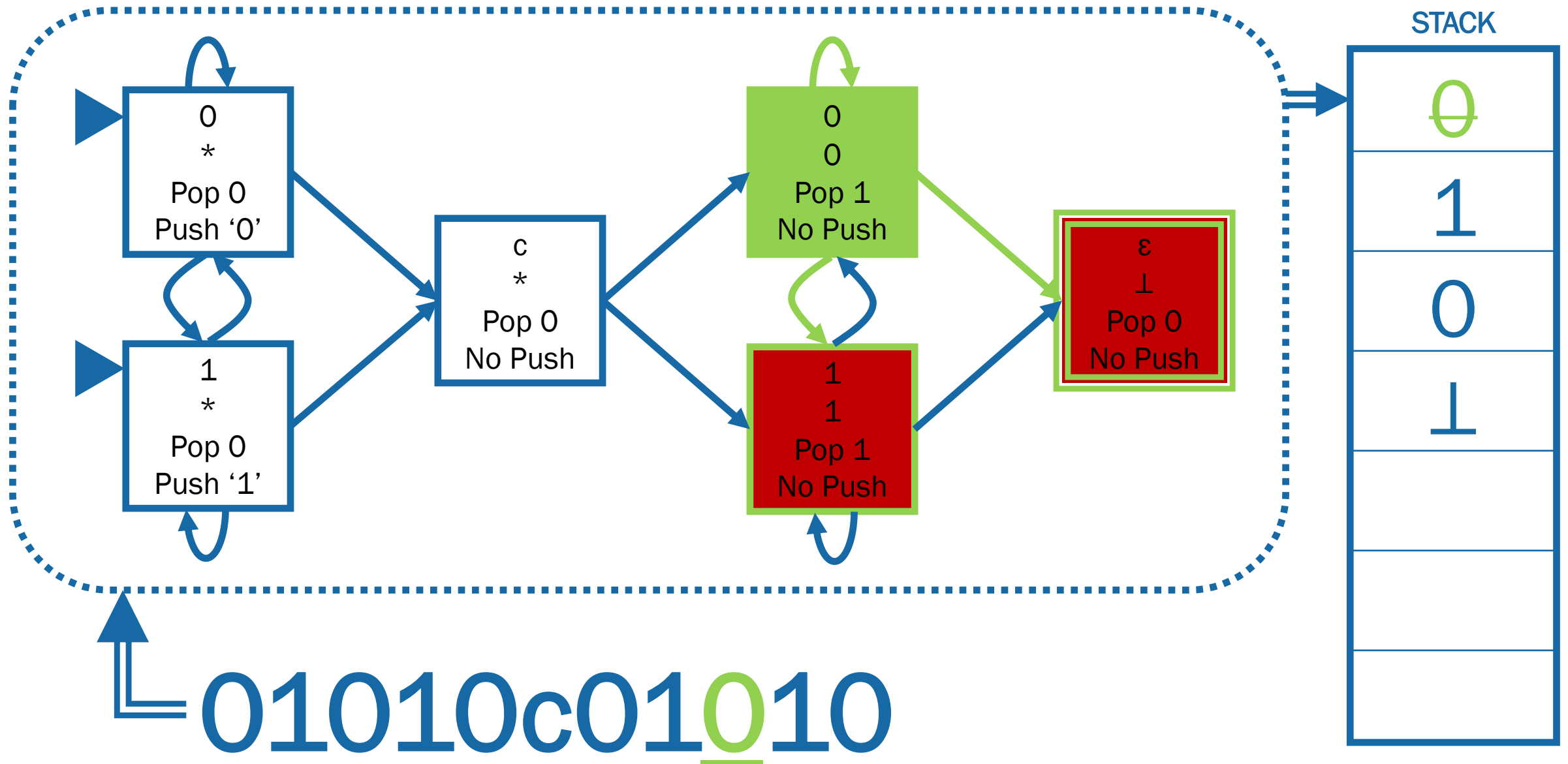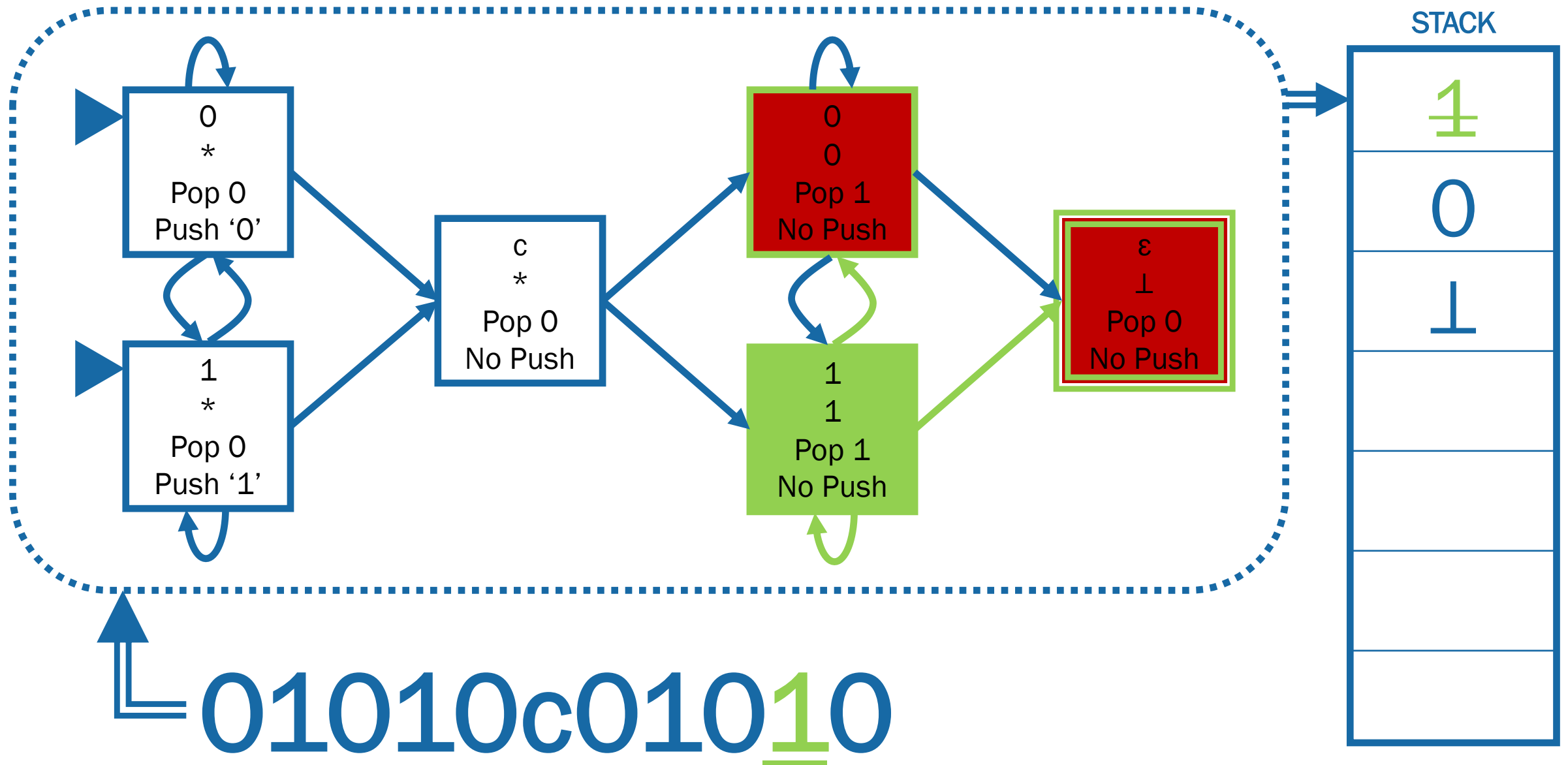


STACK

01010c01010

34

# Recognizing Palindromes with a Middle Character

# Recognizing Palindromes with a Middle Character



STACK

01010c01010 ✓

36

# Mapping DPDA Efficiently to Hardware

- ASPEN supports homogeneous DPDA
  - All transitions to a state occur on the same input character, stack comparison, and stack operation
  - Similar in nature to homogeneous NFAs

- Equal expressive power as standard DPDA

- State increase is quadratic in the worst case with a fixed alphabet

- Allows for efficient mapping to hardware resources
  - Transitions decoupled from input/stack matches

# Five Steps of DPDA Execution Per Cycle



STACK

01010c01010

# Five Steps of DPDA Execution Per Cycle

1. Input Match$^\varepsilon$
2. Stack Match
3. Action Lookup
4. Stack Update
5. State Transition



STACK

01010c01010

$1^\varepsilon$

39

# Five Steps of DPDA Execution Per Cycle

1. Input Match$^\varepsilon$
2. Stack Match
3. Action Lookup
4. Stack Update
5. State Transition



STACK

0
*
Pop 0
Push '0'

1
*
Pop 0
Push '1'

0
*
Pop 0
No Push

0
0
Pop 1
No Push

1
1
Pop 1
No Push

0
⊥
Pop 0
No Push

⊥

**2**

**1**$^\varepsilon$

01010c01010

# Five Steps of DPDA Execution Per Cycle

1. Input Match$^{\varepsilon}$
2. Stack Match
3. Action Lookup
4. Stack Update
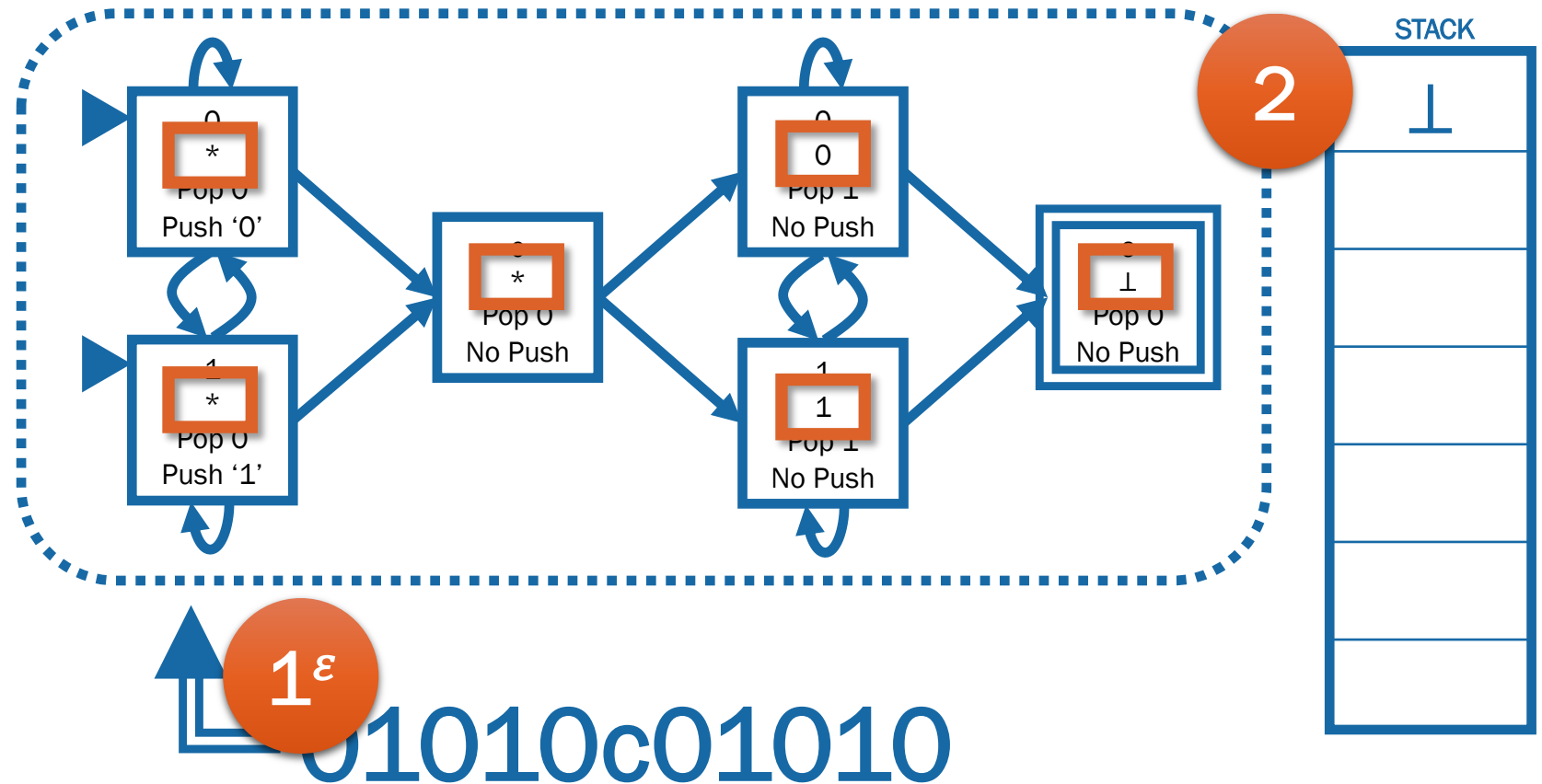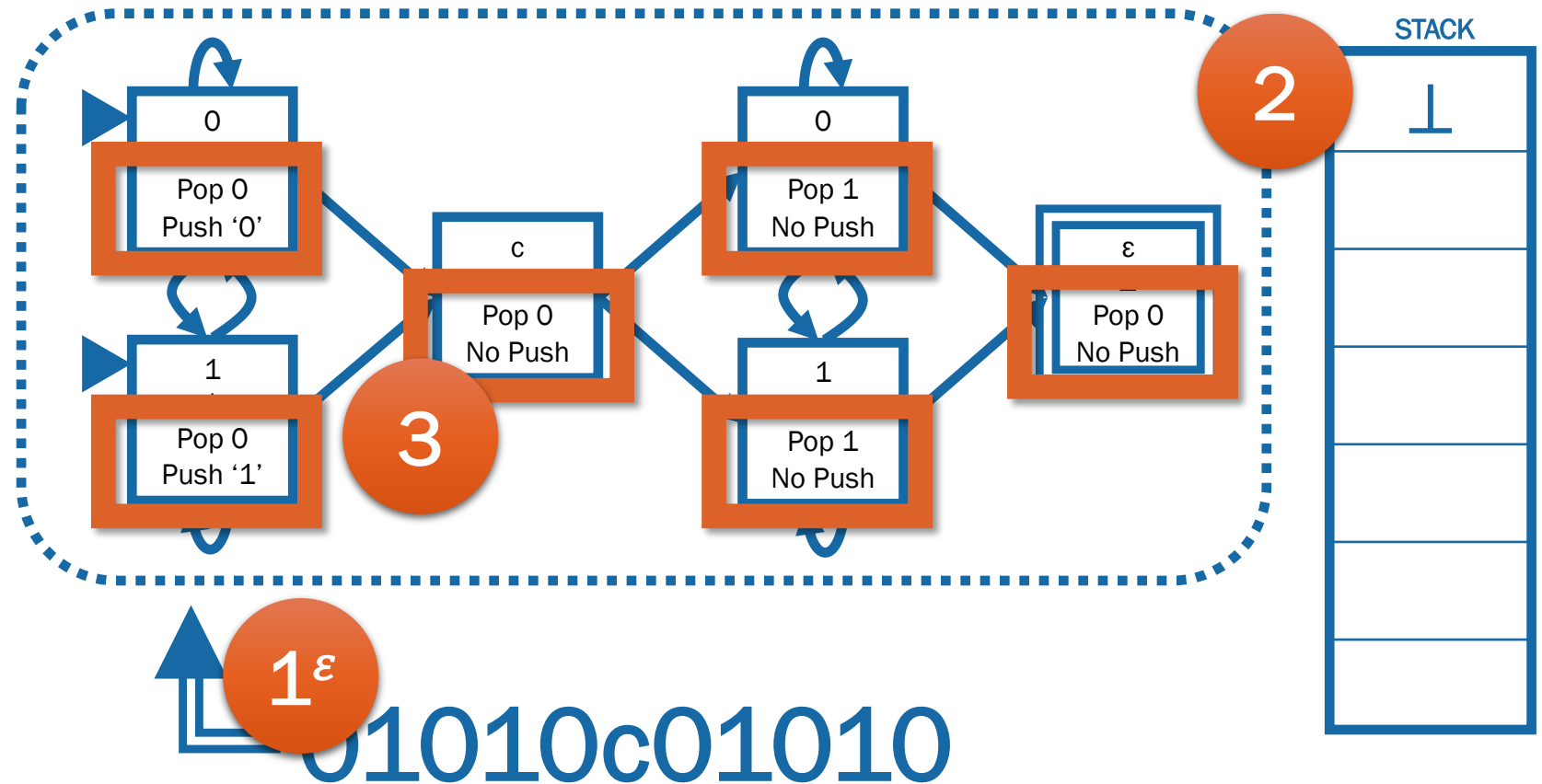5. State Transition

# Five Steps of DPDA Execution Per Cycle

1. Input Match$^{\varepsilon}$
2. Stack Match
3. Action Lookup
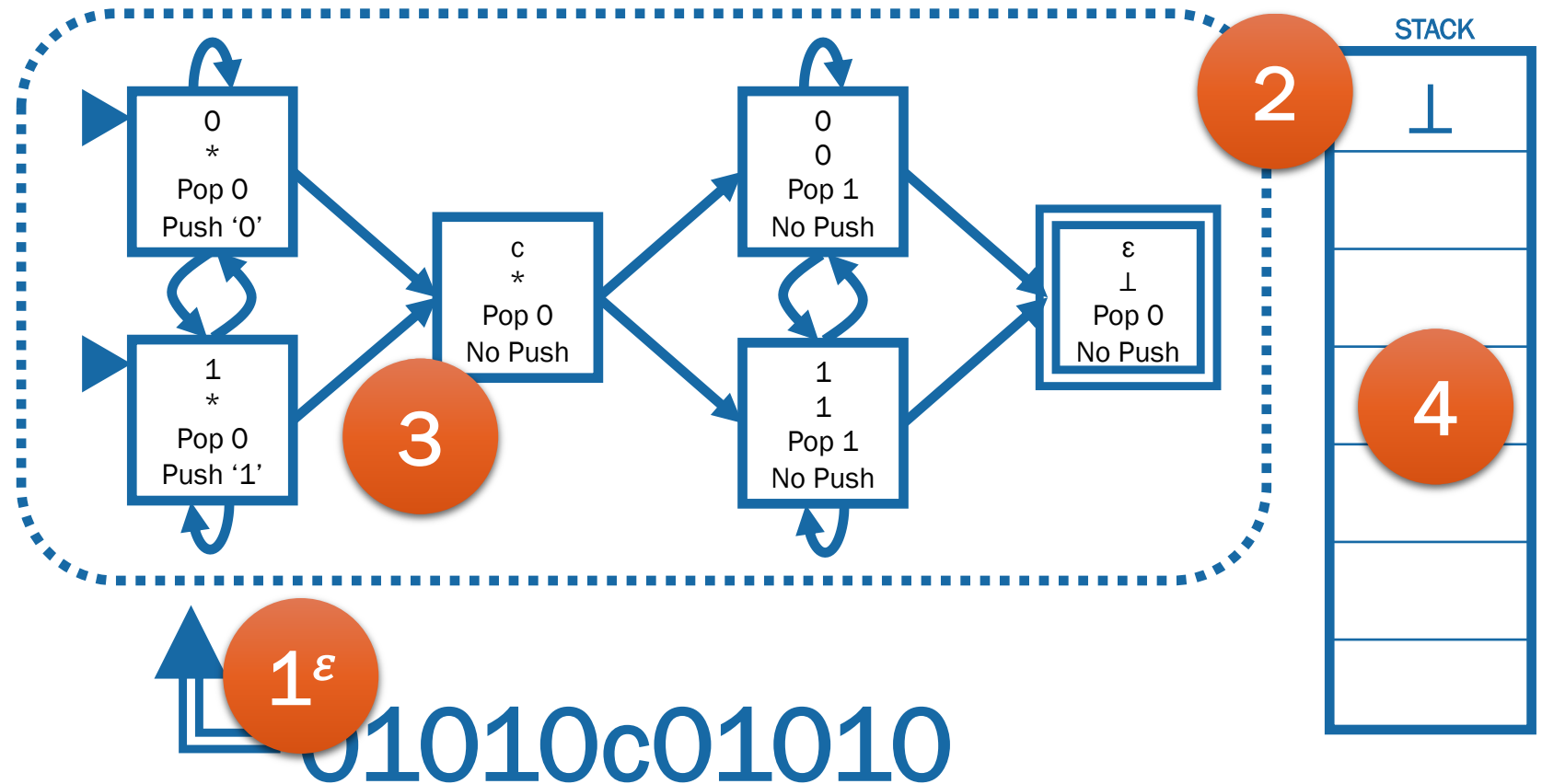4. Stack Update
5. State Transition

# Five Steps of DPDA Execution Per Cycle

1. Input Match$^\varepsilon$
2. Stack Match
3. Action Lookup
4. Stack Update
5. State Transition



01010c01010

# Implementing ASPEN in LLC

- ASPEN repurposes LLC slices for pushdown automata computation
- Location in LLC supports tighter coupling with CPU operations than dedicated accelerator
  - PDA often part of a larger workflow
  - ASPEN similar to auxiliary functional unit in CPU (similar to FPU or vector unit)
- SRAM arrays in LLC already support necessary operations for DPDA execution

# Where is ASPEN?



- ASPEN uses 2 arrays per bank
- 240 states per bank
- Full connectivity within bank
- Global switch and stack in CBOX for large DPDA

# Stack Match in SRAM

- Check **all states** against top of stack
  - One column of SRAM/state
  - Input TOS as row address
  - "1": match; "0": no match

- **Intersect** with currently active states



0                                    255

0

Row Decoder

← Top of Stack (TOS)

255

4:1 column mux

Active State Vector

Active States Matching Stack

# Stack Match in SRAM

- Check **all states** against top of stack
  - One column of SRAM/state
  - Input TOS as row address
  - "1": match; "0": no match

- **Intersect** with currently active states



0                    255

Row Decoder ← Top of Stack (TOS)

4:1 column mux

Active State Vector → Active States Matching Stack
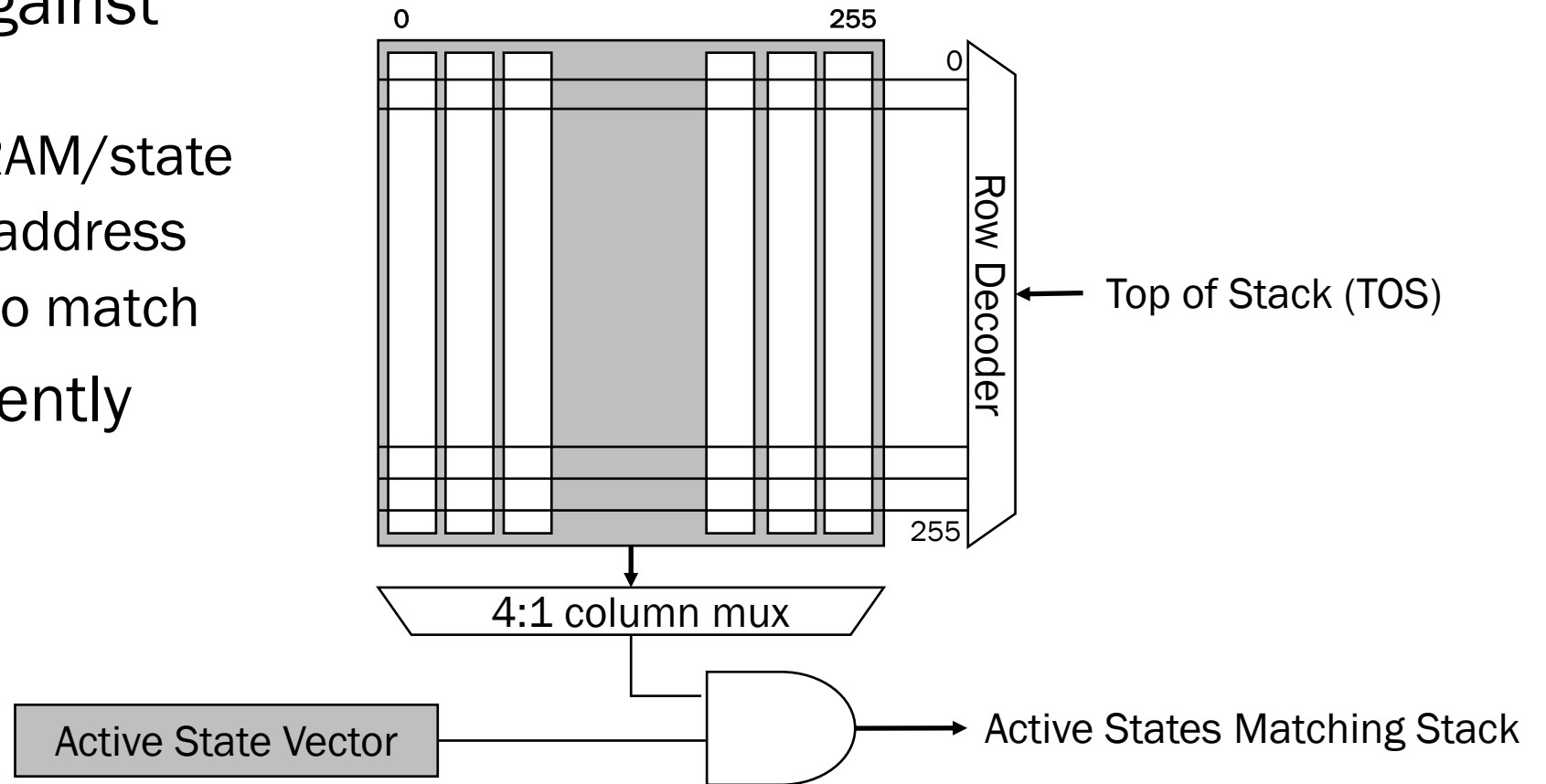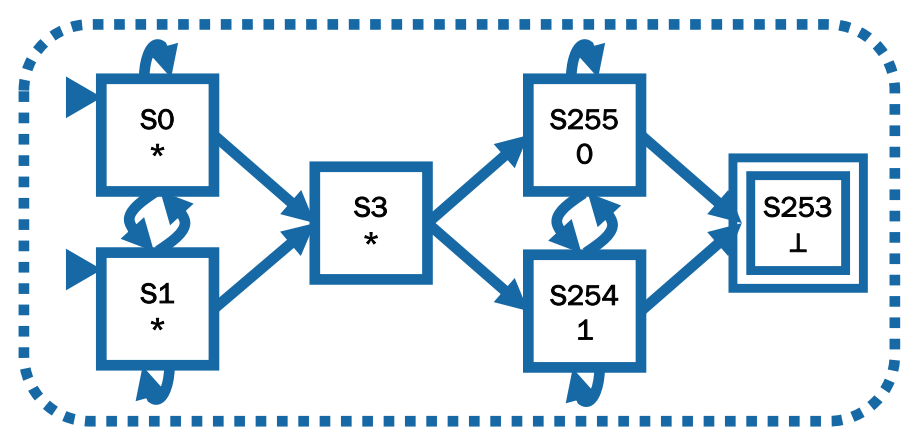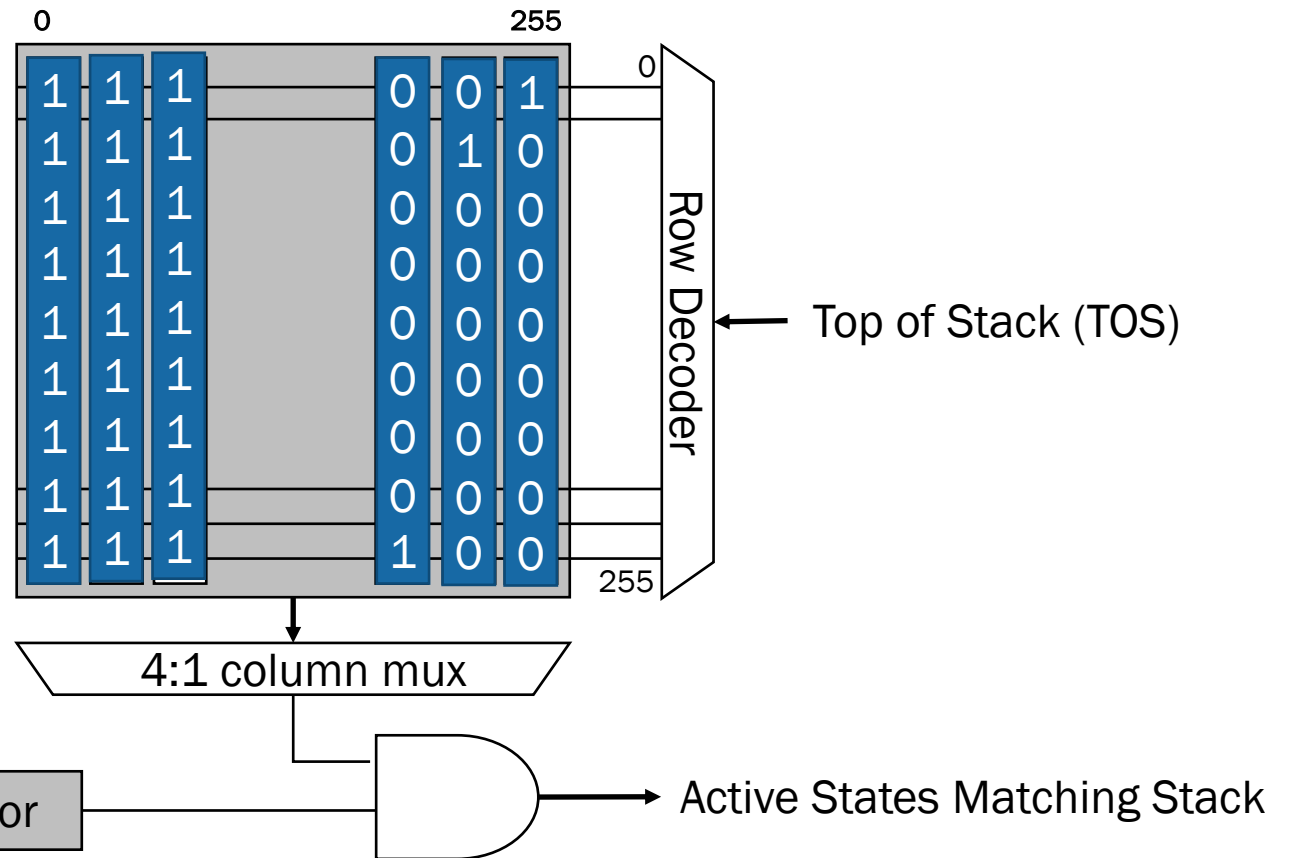
S0 *
S1 *
S3 *
S255 0
S254 1
S253 ⊥

# Stack Match in SRAM



- Check all states against top of stack
  - One column of SRAM/state
  - Input TOS as row address
  - "1": match; "0": no match

- Intersect with currently active states

# ASPEN Datapath — 240 States per Two SRAM Arrays

1. Input Match $\varepsilon$
2. Stack Match
3. Action Lookup
4. Stack Update
5. State Transition

# ASPEN Datapath — 240 States per Two SRAM Arrays



1. Input Match $\varepsilon$
2. Stack Match
3. Action Lookup
4. Stack Update
5. State Transition

# ASPEN Datapath — 240 States per Two SRAM Arrays

1. Input Match$^\varepsilon$
2. Stack Match
3. Action Lookup
4. Stack Update
5. State Transition

# ASPEN Datapath — 240 States per Two SRAM Arrays

1. Input Match$^{\varepsilon}$
2. Stack Match
3. Action Lookup
4. Stack Update
5. State Transition



SRAM Array0

SRAM Array1

0    255

0    255

Stack Pointer

Row Decoder

8-bit Input

1. Input Matching — One Column per State

3. Stack Actions — One Row per State

2. Stack Matching — One Column per State

Row Decoder

0

255

EN

0

255

4:1 column mux

240b    8b    8b

IM Vector

Push Sym.    Pop #

4:1 column mux

240b    8b

SM Vector

TOS +1    TOS

Local TOS

256b

256b

256b

Active State Vector
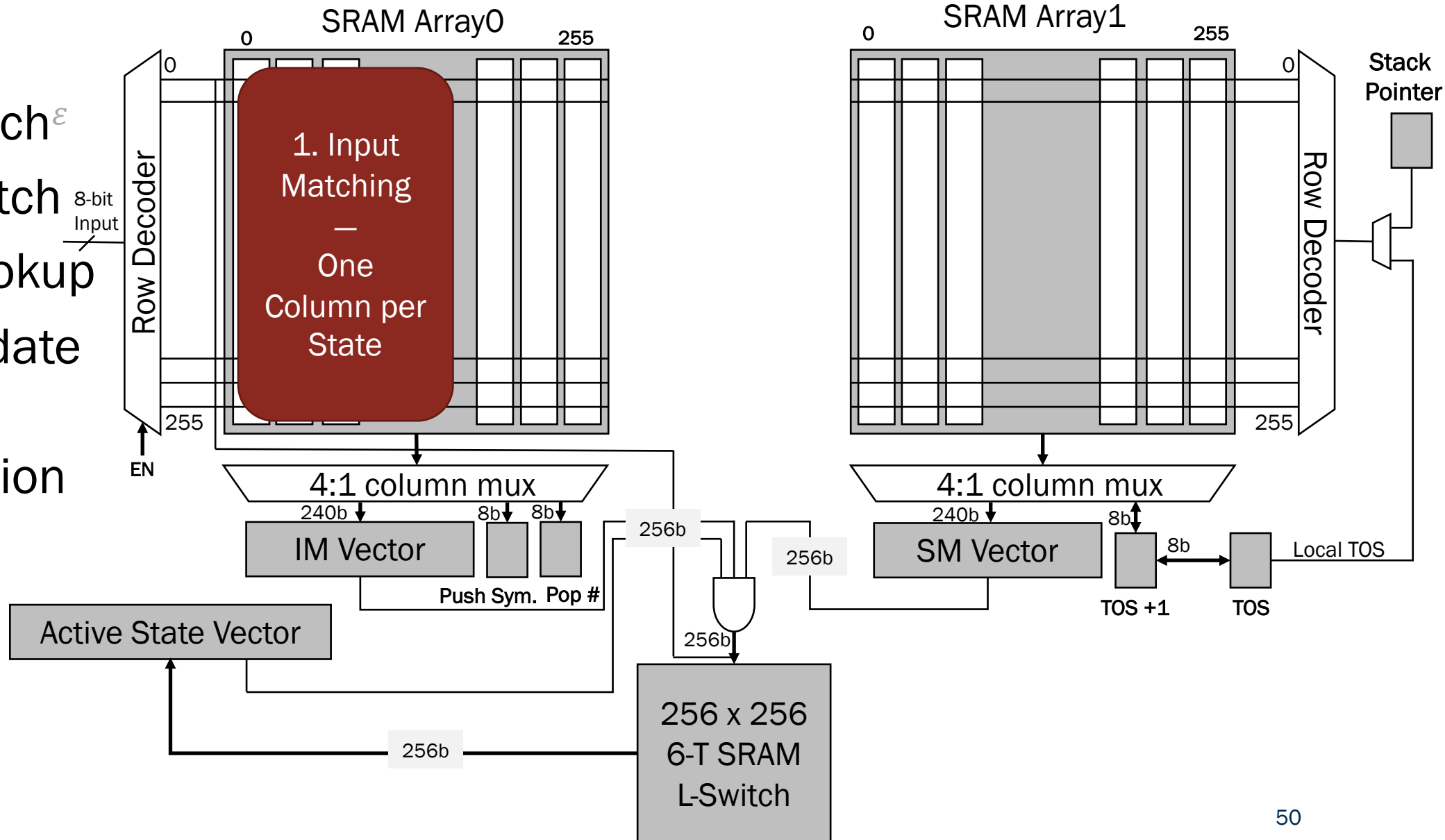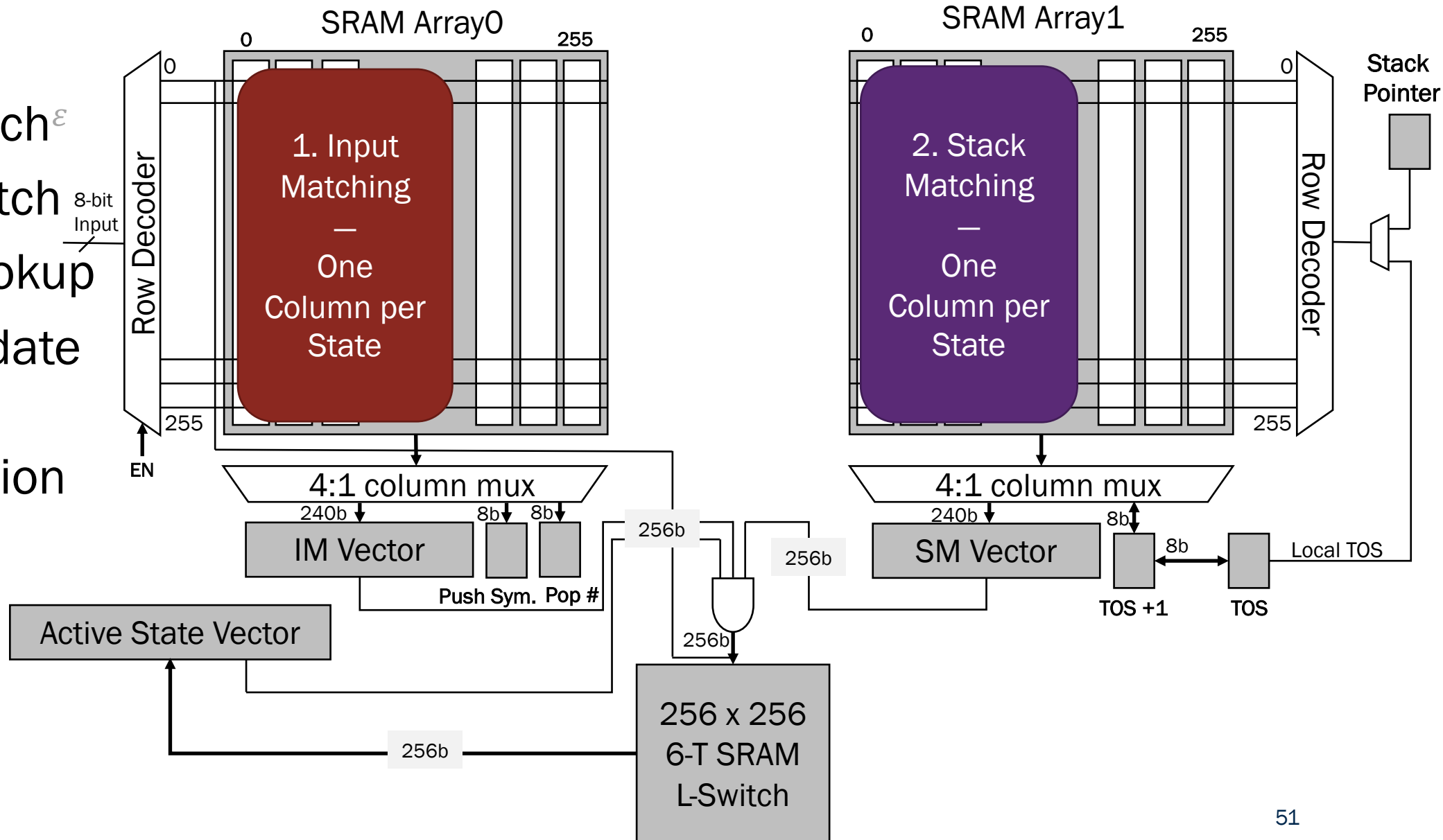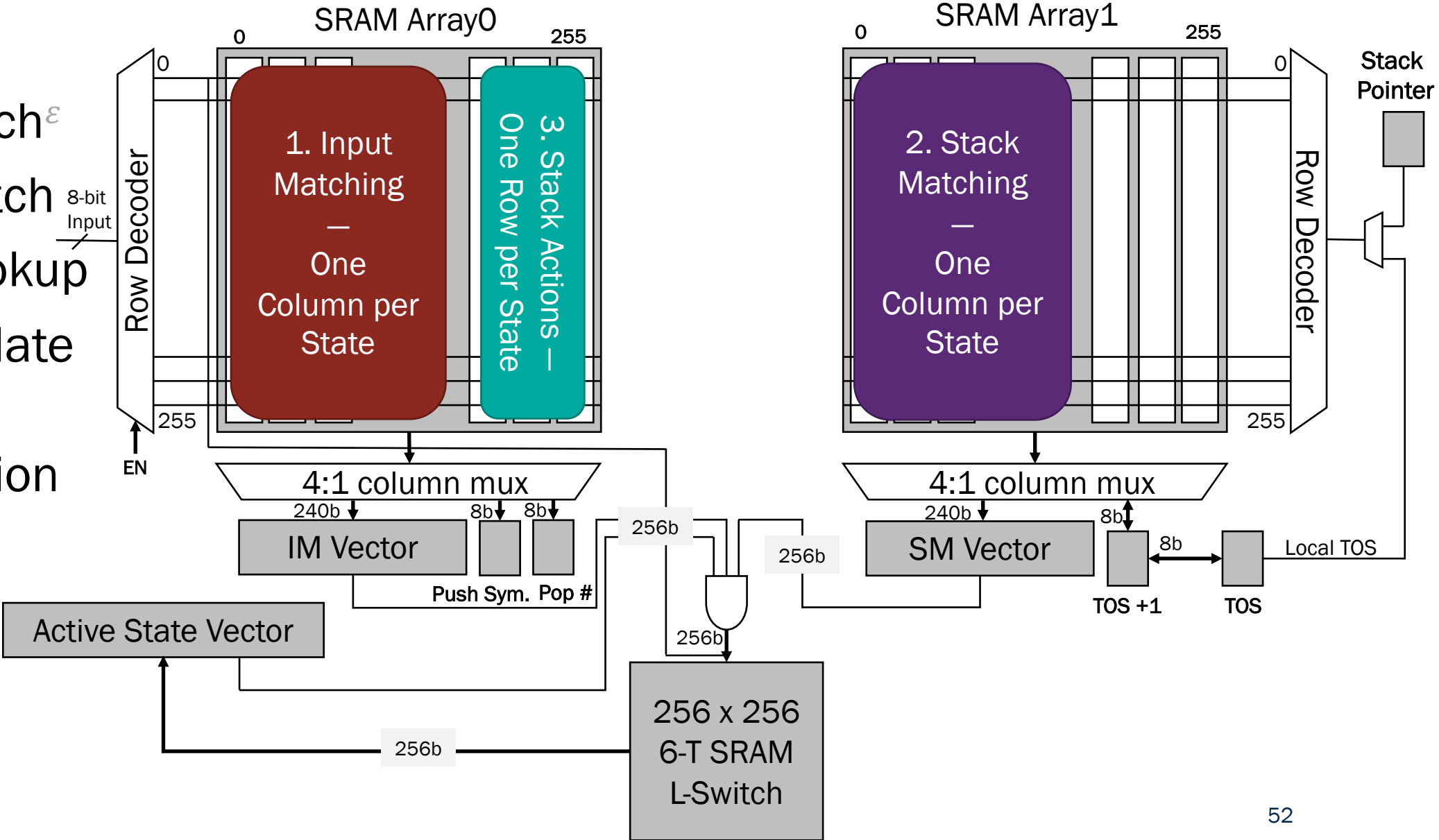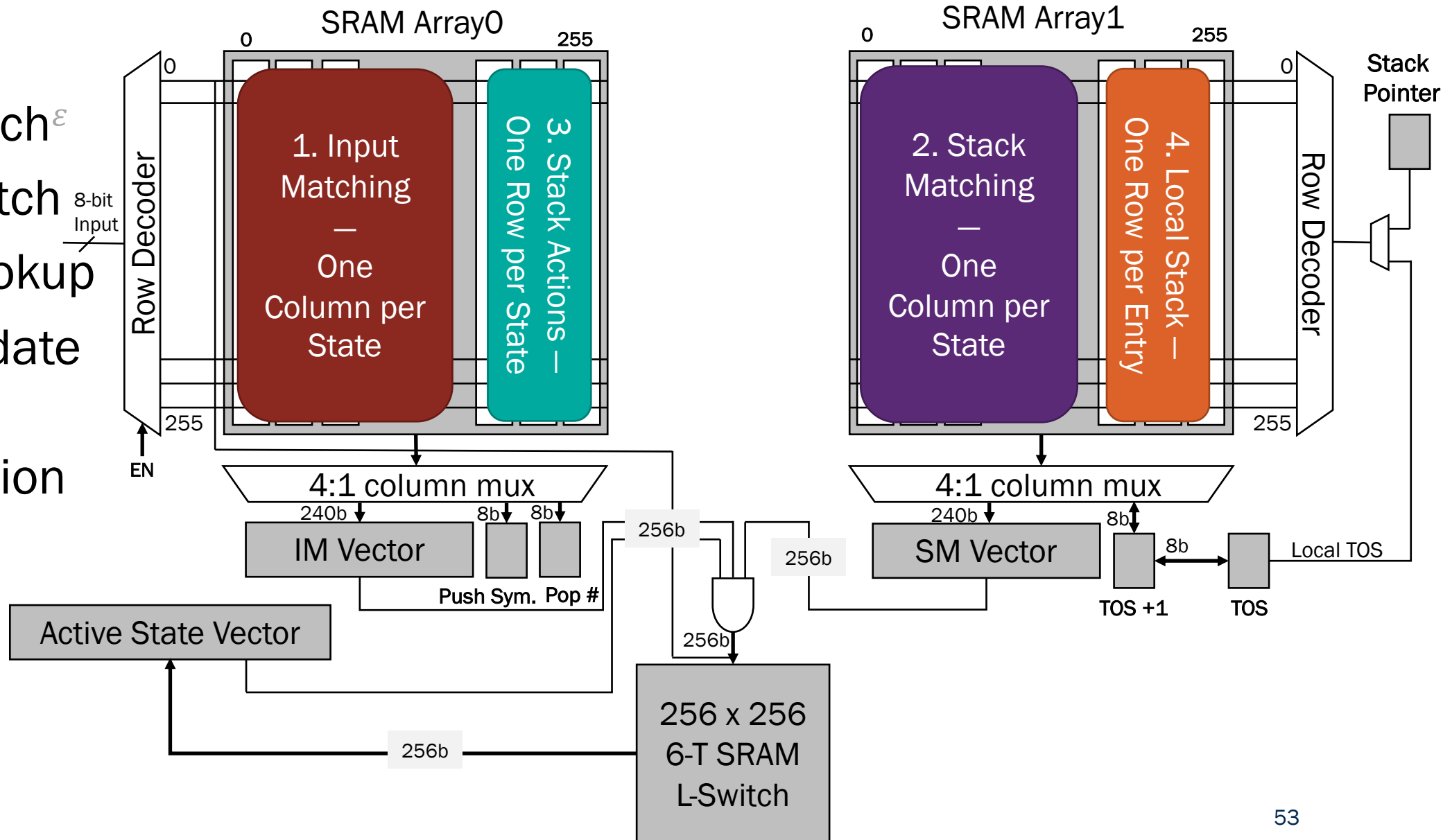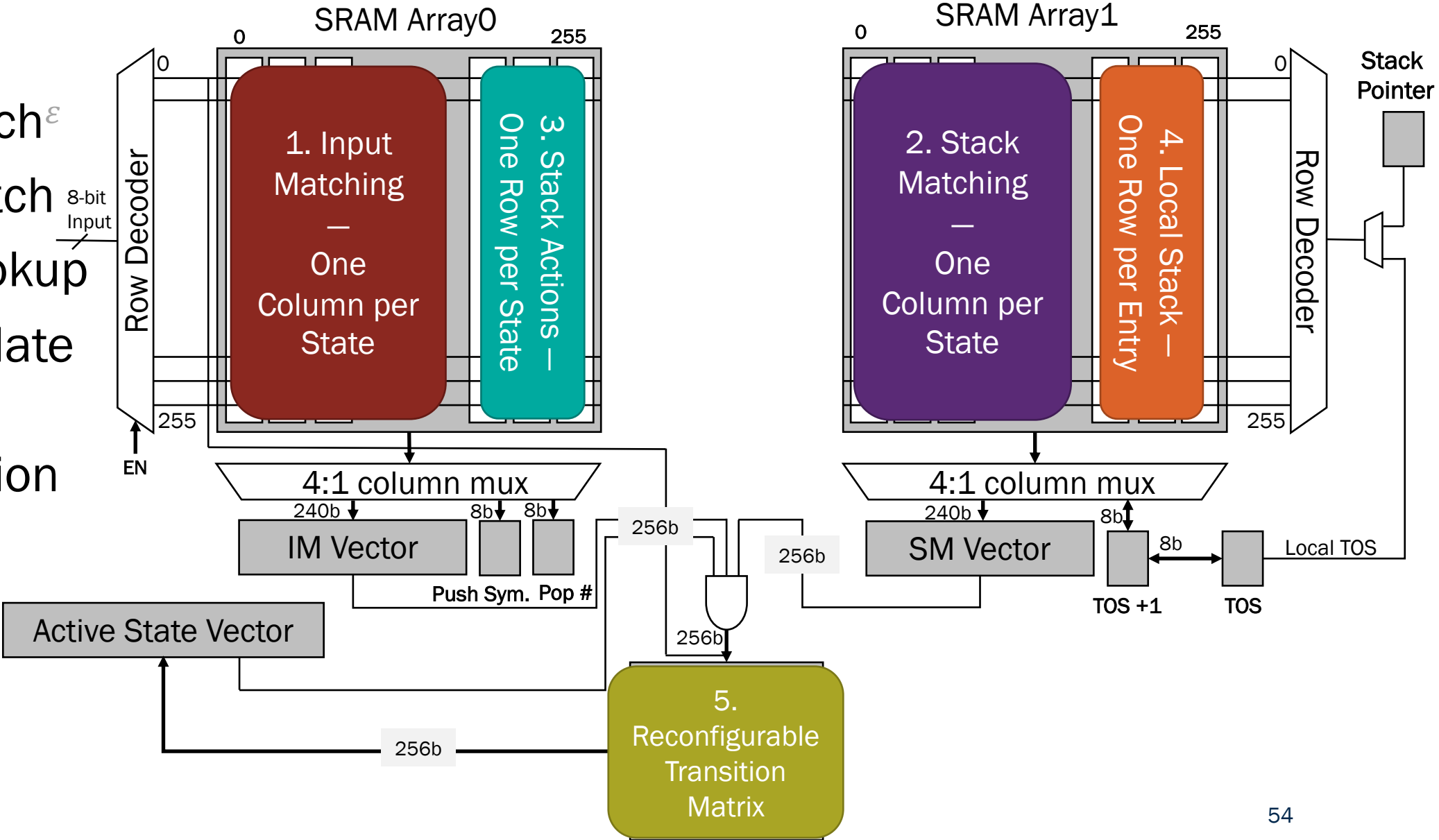
256b

256 x 256 6-T SRAM L-Switch

# ASPEN Datapath — 240 States per Two SRAM Arrays

1. Input Match$^\varepsilon$
2. Stack Match
3. Action Lookup
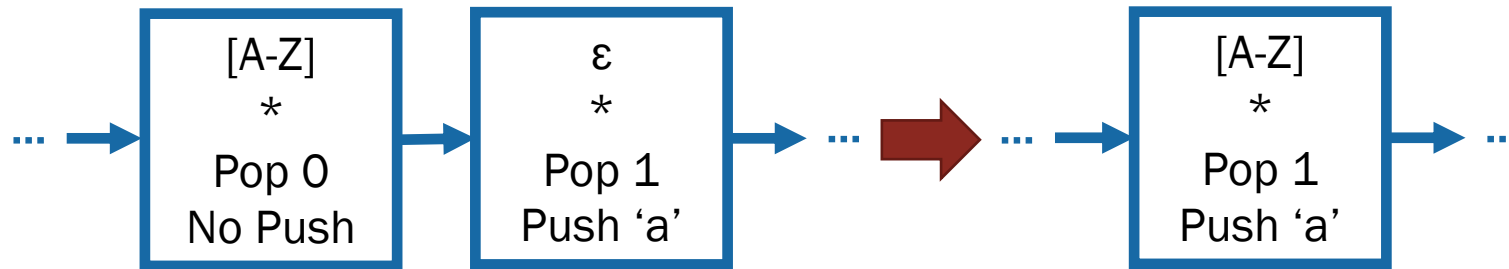4. Stack Update
5. State Transition

# ASPEN Datapath — 240 States per Two SRAM Arrays

1. Input Match$^\varepsilon$
2. Stack Match
3. Action Lookup
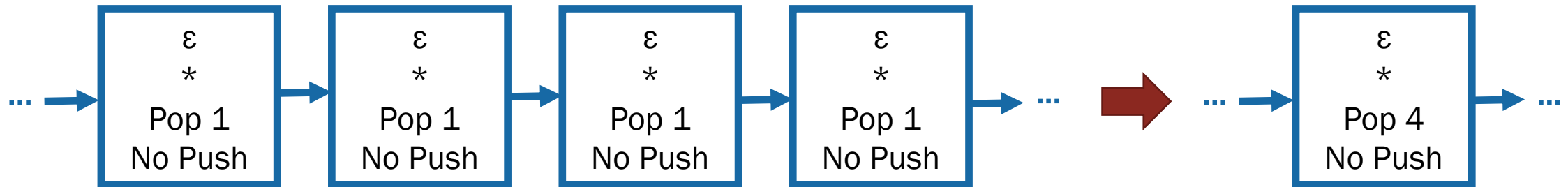4. Stack Update
5. State Transition

# Optimizations

## Epsilon Merging

```
        ┌─────────┐      ┌─────────┐              ┌─────────┐
        │  [A-Z]  │      │    ε    │              │  [A-Z]  │
 ...──▶ │    *    │ ──▶  │    *    │ ──▶ ...  ⟹  ...──▶│    *    │ ──▶ ...
        │  Pop 0  │      │  Pop 1  │              │  Pop 1  │
        │ No Push │      │ Push 'a'│              │ Push 'a'│
        └─────────┘      └─────────┘              └─────────┘
```

## Multipop

```
     ┌─────────┐    ┌─────────┐    ┌─────────┐    ┌─────────┐              ┌─────────┐
     │    ε    │    │    ε    │    │    ε    │    │    ε    │              │    ε    │
...─▶│    *    │─▶  │    *    │─▶  │    *    │─▶  │    *    │─▶ ... ⟹ ...─▶│    *    │─▶ ...
     │  Pop 1  │    │  Pop 1  │    │  Pop 1  │    │  Pop 1  │              │  Pop 4  │
     │ No Push │    │ No Push │    │ No Push │    │ No Push │              │ No Push │
     └─────────┘    └─────────┘    └─────────┘    └─────────┘              └─────────┘
```

**Goal:** Reduce the number of stalls while processing input

- Average of 65% reduction in epsilon states

# Evaluation: Two Real-World Applications

- XML Parsing
  - Common to many data analyses (needed to read input data)
  - Step in a larger pipeline: Tokenization, Parsing, Validation, DOM construction
  - Pipelined with Cache Automaton for tokenization
  - Single Large DPDA with Global Stack

- Frequent Subtree Mining
  - Task of identifying subtrees occurring above a threshold frequency in a corpus of trees
  - Common in recommendation systems, packet routing, NLP, etc.
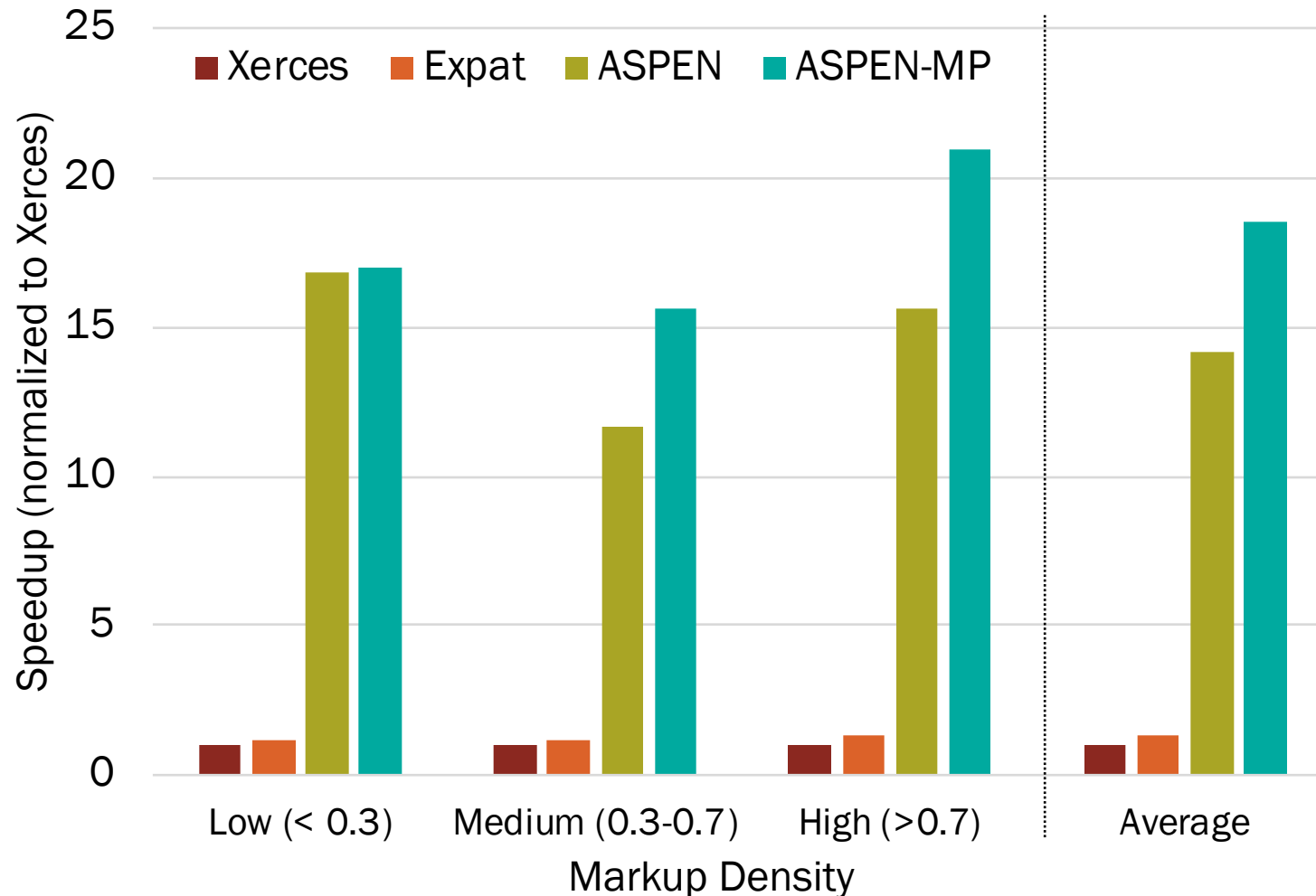  - Many Small DPDA with Local Stacks

# Experimental Methodology

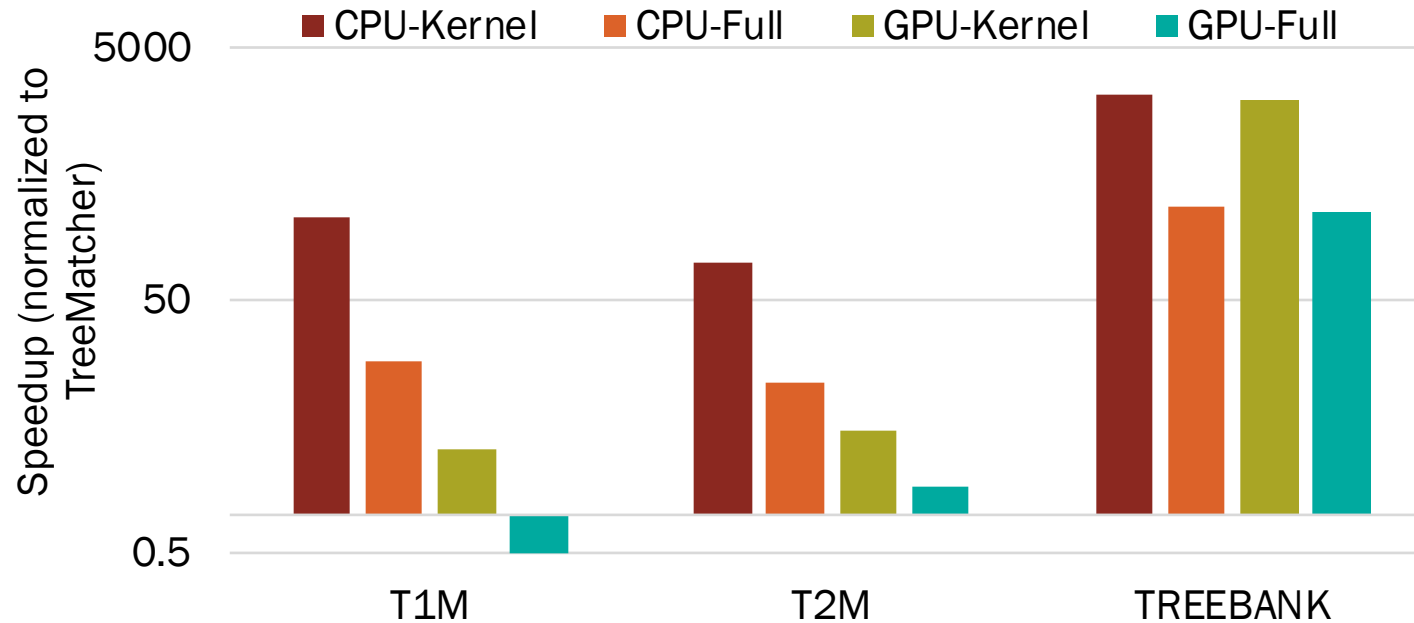| Component | Max Frequency | Operating Frequency |
|---|---|---|
| ASPEN | 880 MHz | 850 MHz |
| Cache Automaton | 4 GHz | 3.4 GHz |

- Baseline Evaluation
  - **CPU:** 2.6 GHz dual-socket Intel Xeon E5-2697-v3 (28 cores total)
  - **GPU:** NVIDIA TITAN Xp
- **Performance and Power:** PAPI, Intel RAPL, NVIDIA *nvprof*
- **ASPEN Simulation:**
  - METIS graph partitioning framework
  - VASim modified for cycle-accurate DPDA simulation

# XML Parsing: Parabix, Ximpleware, UW XML



- ASPEN is 13-18x faster (on average) than popular CPU Parsers

- Performance did not vary significantly with complexity of XML

- Optimizations and tokenization hide $\varepsilon$-stalls
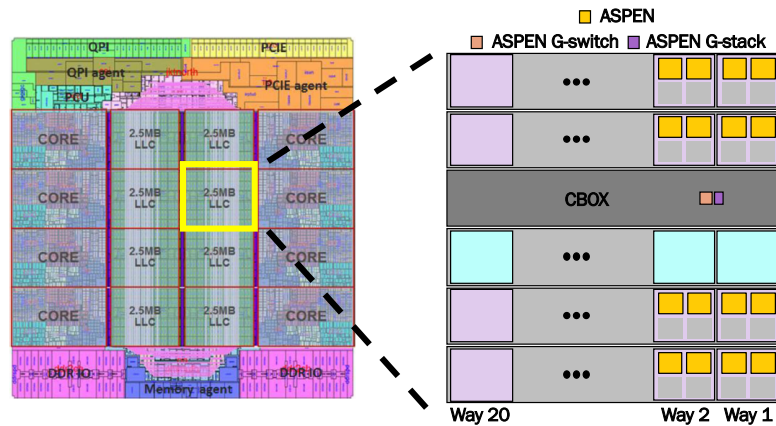
# Frequent Subtree Mining



Legend: CPU-Kernel | CPU-Full | GPU-Kernel | GPU-Full

Y-axis: Speedup (normalized to TreeMatcher) — 5000, 50, 0.5

X-axis: T1M, T2M, TREEBANK

| Dataset | Automata Alphabets | Stack Alphabets | Stack Size |
|---------|--------------------|-----------------|------------|
| T1M | 16 | 17 | 29 |
| T2M | 38 | 39 | 49 |
| TREEBANK | 100 | 101 | 110 |

- ASPEN is (on average) 67x faster than CPUs 6x faster than GPUs for end-to-end application

- Performance on ASPEN is independent from tree size and complexity

- No $\varepsilon$-transitions

# Conclusion

# Conclusion



ASPEN: Processor for DPDA Acceleration

# Conclusion



$S \rightarrow Exp \dashv$
$Exp \rightarrow Term + Exp$
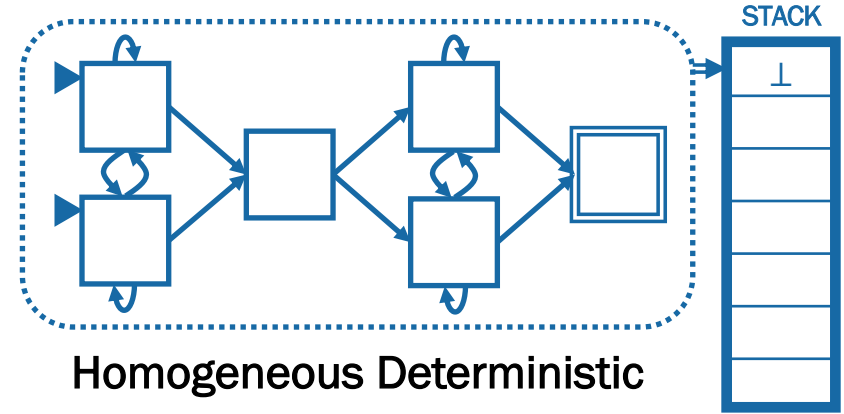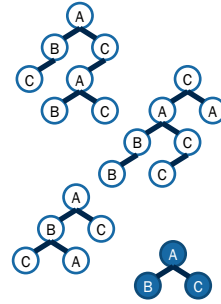$| Term$
$Term \rightarrow$ `int` $* Term$
$| ( Exp )$
$|$ `int`

Supports Processing of
Recursively-Nested and Tree-Structured Data



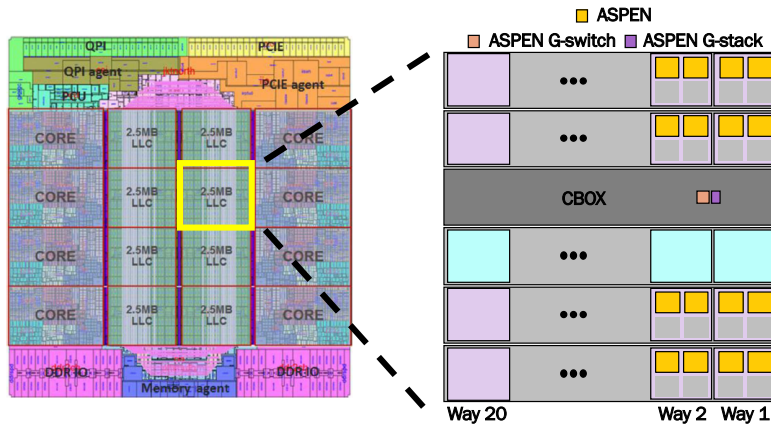ASPEN: Processor for DPDA Acceleration

# Conclusion



$S \rightarrow Exp \dashv$
$Exp \rightarrow Term + Exp$
$\quad | \ Term$
$Term \rightarrow \texttt{int} * Term$
$\quad | \ ( \ Exp \ )$
$\quad | \ \texttt{int}$

Supports Processing of
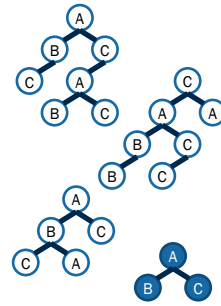Recursively-Nested and Tree-Structured Data

Homogeneous Deterministic
Pushdown Automaton
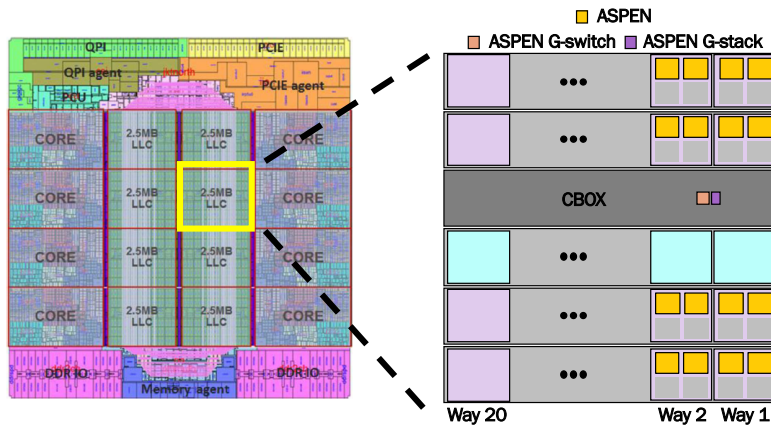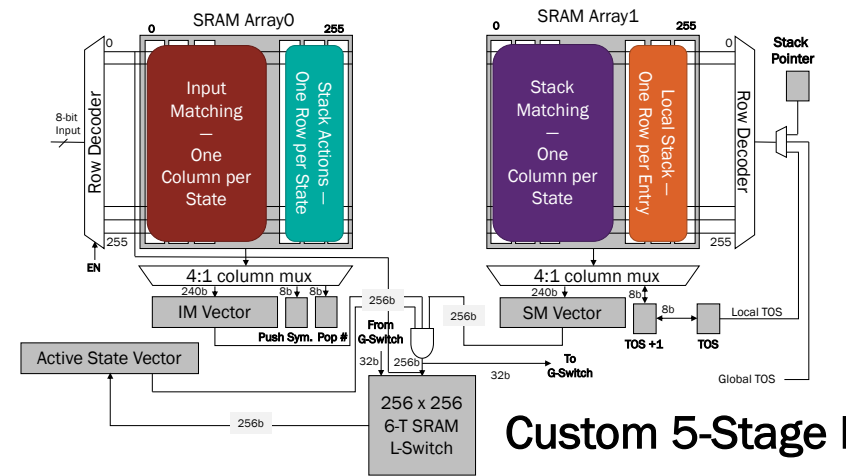
ASPEN: Processor for DPDA Acceleration

# Conclusion



$S \rightarrow Exp \dashv$
$Exp \rightarrow Term + Exp$
$\quad | Term$
$Term \rightarrow \texttt{int} * Term$
$\quad | ( Exp )$
$\quad | \texttt{int}$

Supports Processing of
Recursively-Nested and Tree-Structured Data

Homogeneous Deterministic
Pushdown Automaton

STACK

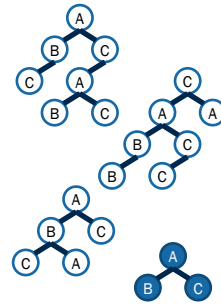ASPEN: Processor for DPDA Acceleration

Custom 5-Stage Datapath

64

# Conclusion

$S \rightarrow Exp \dashv$
$Exp \rightarrow Term + Exp$
$| Term$
$Term \rightarrow \texttt{int} * Term$
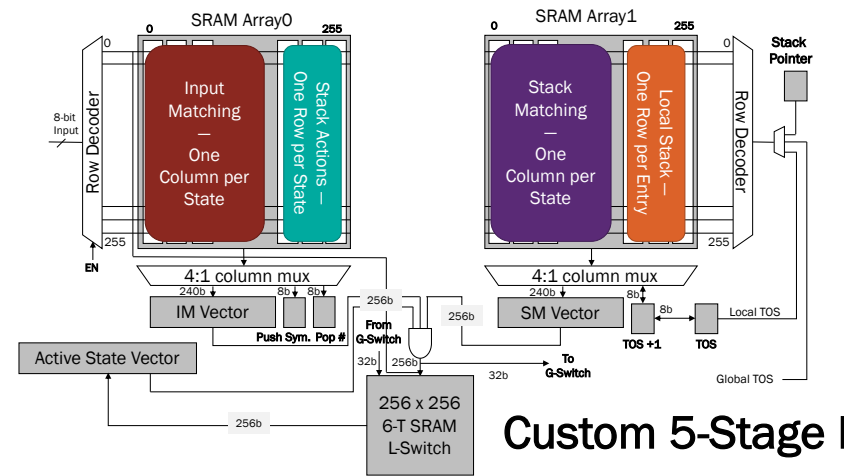$| ( Exp )$
$| \texttt{int}$

Supports Processing of
Recursively-Nested and Tree-Structured Data

Homogeneous Deterministic
Pushdown Automaton

STACK

ASPEN: Processor for DPDA Acceleration

Custom 5-Stage Datapath