**DVLabs**

SCREAM

Rob King

Introduction
The Problem
The Base64
Algorithm
Regular Expressions

The Algorithm
Ways of Solving the
Problem
Encoding Operations

Performance
Expression
Optimization
Performance
Analysis

Implementation
and Usage
Common Use Cases
Caveats

Summary

# Static Compilation of Regular Expressions for Analysis and Modification

Rob King

DVLabs
TippingPoint Technologies

February 18, 2010

# OUTLINE

SCREAM

Rob King

Introduction

The Problem
The Base64
Algorithm
Regular Expressions

The Algorithm

Ways of Solving the
Problem
Encoding Operations

Performance

Expression
Optimization
Performance
Analysis

Implementation
and Usage

Common Use Cases
Caveats

Summary

**DVLabs**

THE PROBLEM
In which we discover purpose of the whole thing...

- This talk is about inspecting streams of data for interesting patterns, even when that stream of data has been encoded.

- We focus on the Base64 encoding scheme, and discuss a tool that can be used when dealing with Base64.

- However, most portions of the algorithm are applicable to other position-dependent bitwise block encodings (and, potentially, self-synchronizing encodings).

**DVLabs**

THE PROBLEM
In which we discover purpose of the whole thing...

- This talk is about inspecting streams of data for interesting patterns, even when that stream of data has been encoded.

- We focus on the Base64 encoding scheme, and discuss a tool that can be used when dealing with Base64.

- However, most portions of the algorithm are applicable to other position-dependent bitwise block encodings (and, potentially, self-synchronizing encodings).

**DVLabs**

THE PROBLEM
In which we discover purpose of the whole thing...

- This talk is about inspecting streams of data for interesting patterns, even when that stream of data has been encoded.
- We focus on the Base64 encoding scheme, and discuss a tool that can be used when dealing with Base64.
- However, most portions of the algorithm are applicable to other position-dependent bitwise block encodings (and, potentially, self-synchronizing encodings).

# THE PROBLEM
In which we discover purpose of the whole thing...

- This talk is about inspecting streams of data for interesting patterns, even when that stream of data has been encoded.
- We focus on the Base64 encoding scheme, and discuss a tool that can be used when dealing with Base64.
- However, most portions of the algorithm are applicable to other position-dependent bitwise block encodings (and, potentially, self-synchronizing encodings).

When looking for patterns in streams of data, several things must be kept in mind:

- There is no "luxury of time".
- Context is limited.
- Resources are limited.

When looking for patterns in streams of data, several things
must be kept in mind:

- There is no "luxury of time".
- Context is limited.
- Resources are limited.

When looking for patterns in streams of data, several things
must be kept in mind:

- There is no "luxury of time".
- Context is limited.
- Resources are limited.

- Sometimes, the streams we're inspecting will be encoded.
- This means that we're going to have to be (more!) clever when looking for patterns in these streams.

There are several general strategies for dealing with encoded streams.

- The easiest thing to do is simply pretend the stream is not encoded at all.
- Advantages:
  - We're already done.
- Disadvantages:
  - We're essentially admitting defeat.
  - Whatever we were looking for is not going to be found.
  - We're stil burdening our analysis engine with lots of data with which we can do nothing.

- The easiest thing to do is simply pretend the stream is not encoded at all.
- Advantages:
  - We're already done.
- Disadvantages:
  - We're essentially admitting defeat.
  - Whatever we were looking for is not going to be found.
  - We're stil burdening our analysis engine with lots of data with which we can do nothing.

- The easiest thing to do is simply pretend the stream is not encoded at all.
- Advantages:
  - We're already done.
- Disadvantages:
  - We're essentially admitting defeat.
  - Whatever we were looking for is not going to be found.
  - We're stil burdening our analysis engine with lots of data with which we can do nothing.

- Rather than pretending that the data is not encoded, we could go one step further, and detect that the data is encoded.
- Once we've detected that the data is encoded, we can simply drop the stream.
- Advantages:
  - Stops burdening with inspection engine with data we know we can't inspect.
- Disadvantages:
  - If we "fail open" and allow encoded data to pass without inspection, we just gave anyone who wants to bypass our inspection a "get out of jail free" card.
  - If we "fail closed" and block encoded data, we just blocked all legitimate uses of that encoding scheme.

- Rather than pretending that the data is not encoded, we could go one step further, and detect that the data is encoded.
- Once we've detected that the data is encoded, we can simply drop the stream.
- Advantages:
  - Stops burdening with inspection engine with data we know we can't inspect.
- Disadvantages:
  - If we "fail open" and allow encoded data to pass without inspection, we just gave anyone who wants to bypass our inspection a "get out of jail free" card.
  - If we "fail closed" and block encoded data, we just blocked all legitimate uses of that encoding scheme.

- Rather than pretending that the data is not encoded, we could go one step further, and detect that the data is encoded.
- Once we've detected that the data is encoded, we can simply drop the stream.
- Advantages:
  - Stops burdening with inspection engine with data we know we can't inspect.
- Disadvantages:
  - If we "fail open" and allow encoded data to pass without inspection, we just gave anyone who wants to bypass our inspection a "get out of jail free" card.
  - If we "fail closed" and block encoded data, we just blocked all legitimate uses of that encoding scheme.

- Rather than pretending that the data is not encoded, we could go one step further, and detect that the data is encoded.
- Once we've detected that the data is encoded, we can simply drop the stream.
- Advantages:
  - Stops burdening with inspection engine with data we know we can't inspect.
- Disadvantages:
  - If we "fail open" and allow encoded data to pass without inspection, we just gave anyone who wants to bypass our inspection a "get out of jail free" card.
  - If we "fail closed" and block encoded data, we just blocked all legitimate uses of that encoding scheme.

- We could buffer the entirety of an encoded stream.
- Once we've buffered the whole stream, we can decode it, inspect it, reencode it, and send it on its way.
- Advantages:
  - We get the complete power of our inspection engine.
- Disadvantages:
  - Latency becomes unbounded.
  - Resource usage becomes unbounded.
  - We have to modify the engine for every encoding we need to inspect.
- These advantages and disadvantages apply roughly equally to any streaming decoder.

- We could buffer the entirety of an encoded stream.
- Once we've buffered the whole stream, we can decode it, inspect it, reencode it, and send it on its way.
- Advantages:
  - We get the complete power of our inspection engine.
- Disadvantages:
  - Latency becomes unbounded.
  - Resource usage becomes unbounded.
  - We have to modify the engine for every encoding we need to inspect.
- These advantages and disadvantages apply roughly equally to any streaming decoder.

- We could buffer the entirety of an encoded stream.
- Once we've buffered the whole stream, we can decode it, inspect it, reencode it, and send it on its way.
- Advantages:
  - We get the complete power of our inspection engine.
- Disadvantages:
  - Latency becomes unbounded.
  - Resource usage becomes unbounded.
  - We have to modify the engine for every encoding we need to inspect.
- These advantages and disadvantages apply roughly equally to any streaming decoder.

- We could buffer the entirety of an encoded stream.
- Once we've buffered the whole stream, we can decode it, inspect it, reencode it, and send it on its way.
- Advantages:
  - We get the complete power of our inspection engine.
- Disadvantages:
  - Latency becomes unbounded.
  - Resource usage becomes unbounded.
  - We have to modify the engine for every encoding we need to inspect.
- These advantages and disadvantages apply roughly equally to any streaming decoder.

- We could buffer the entirety of an encoded stream.
- Once we've buffered the whole stream, we can decode it, inspect it, reencode it, and send it on its way.
- Advantages:
  - We get the complete power of our inspection engine.
- Disadvantages:
  - Latency becomes unbounded.
  - Resource usage becomes unbounded.
  - We have to modify the engine for every encoding we need to inspect.
- These advantages and disadvantages apply roughly equally to any streaming decoder.

- Instead of decoding the data, we could encode the pattern.
- Advantages:
  - We get (most?) of the power of our inspection engine.
  - There is no performance penalty for decoding.
  - Resource usage and latency are bounded.
  - We need not buffer or store context.
- Disadvantages:
  - A different transform must be written for each scheme.
  - Might not be possible for some schemes or patterns.
  - False positives may become more common.

- Instead of decoding the data, we could encode the pattern.
- Advantages:
  - We get (most?) of the power of our inspection engine.
  - There is no performance penalty for decoding.
  - Resource usage and latency are bounded.
  - We need not buffer or store context.
- Disadvantages:
  - A different transform must be written for each scheme.
  - Might not be possible for some schemes or patterns.
  - False positives may become more common.

- Instead of decoding the data, we could encode the pattern.
- Advantages:
    - We get (most?) of the power of our inspection engine.
    - There is no performance penalty for decoding.
    - Resource usage and latency are bounded.
    - We need not buffer or store context.
- Disadvantages:
    - A different transform must be written for each scheme.
    - Might not be possible for some schemes or patterns.
    - False positives may become more common.

**BASE64**
In which we discover the usefulness of large radices
or "Maybe the Sumerians were right after all..."

**DVLabs**

- Base64 is an Internet standard for encoding arbitrary (usually binary) data using only printable characters common in many character sets.

- Multiple minor variants, but the most common is defined in RFC4648.

- Base64 is used in numerous situations, essentially whenever binary data needs to be encoded in a printable form.

# The most common application of Base64 is in the encoding of attachments to email messages.

```
--Apple-Mail-35-294643828
Content-Disposition: inline;
        filename=stage1.png
Content-Type: image/png;
        x-mac-hide-extension=yes;
        x-unix-mode=0644;
        name="stage1.png"
Content-Transfer-Encoding: base64
```

iVBORw0KGgoAAAANSUhEUgAAAGAAAAIbCAYAAAAU4psTAAAACXBIWXMAAAsTAAALEwEAmpwYAAAg
AEIEQVR4Ae2dCfQUxfHH2yTGnF6JRqNxoJGIkiQbBG/ECFFBBUbzQiKJ4P4RL4yiENI7zofEpp4gFR
8QQUFPEED5AoKhFQPEnAmEMTJRdiYhJYv/rU/nPvv3tb3v3tb3u/tb1nPbuv3tb/n5...
vvAFPX7tqd18zW265pfnyi78cOW+fCZtKbAbyO3nzzTbNiXl/XlxkX9aeOV/mkTV/5+oOTT
V7//6VbPVVluZ73bzhnQ6+s+a73xE 0mi9v9dbbbx2 2nvLJMYfEMVIyb4 ...
```

- Base64 expects input as a series of eight-bit octets.

- Every three octets are grouped together into a collection of 24-bits.

- These 24-bits are then split into four six-bit sextets.

- Each sextet is used as a big-endian index into the zero-based array that is the Base64 alphabet.

| f | | | | | | | | o | | | | | | | | o | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 |
| Z | | | | | | m | | | | | | 9 | | | | | | v | | | | | |

- If the input is not evenly divisible by three, it is right-padded by one or two octets of zeroes, and one or two "=" symbols are appended to the output.

- Base64 expects input as a series of eight-bit octets.
- Every three octets are grouped together into a collection of 24-bits.
- These 24-bits are then split into four six-bit sextets.
- Each sextet is used as a big-endian index into the zero-based array that is the Base64 alphabet.

| f | | | | | | | | o | | | | | | | | o | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 |
| Z | | | | | | m | | | | | | 9 | | | | | | v | | | | | |

- If the input is not evenly divisible by three, it is right-padded by one or two octets of zeroes, and one or two "=" symbols are appended to the output.

- Base64 expects input as a series of eight-bit octets.
- Every three octets are grouped together into a collection of 24-bits.
- These 24-bits are then split into four six-bit sextets.
- Each sextet is used as a big-endian index into the zero-based array that is the Base64 alphabet.

| f | | | | | | | | o | | | | | | | | o | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 |
| Z | | | | | | m | | | | | | 9 | | | | | | v | | | | | |

- If the input is not evenly divisible by three, it is right-padded by one or two octets of zeroes, and one or two "=" symbols are appended to the output.

- Base64 expects input as a series of eight-bit octets.
- Every three octets are grouped together into a collection of 24-bits.
- These 24-bits are then split into four six-bit sextets.
- Each sextet is used as a big-endian index into the zero-based array that is the Base64 alphabet.

| f | | | | | | | | o | | | | | | | | o | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 |
| Z | | | | | | m | | | | | | 9 | | | | | | v | | | | | |

- If the input is not evenly divisible by three, it is right-padded by one or two octets of zeroes, and one or two "=" symbols are appended to the output.

- Base64 expects input as a series of eight-bit octets.
- Every three octets are grouped together into a collection of 24-bits.
- These 24-bits are then split into four six-bit sextets.
- Each sextet is used as a big-endian index into the zero-based array that is the Base64 alphabet.

| f | | | | | | | | o | | | | | | | | o | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 |
| Z | | | | | | m | | | | | | 9 | | | | | | v | | | | | |

- If the input is not evenly divisible by three, it is right-padded by one or two octets of zeroes, and one or two "=" symbols are appended to the output.

- Base64 expects input as a series of eight-bit octets.
- Every three octets are grouped together into a collection of 24-bits.
- These 24-bits are then split into four six-bit sextets.
- Each sextet is used as a big-endian index into the zero-based array that is the Base64 alphabet.

| f | | | | | | | | o | | | | | | | | o | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 |
| Z | | | | | | m | | | | | | 9 | | | | | | v | | | | | |

- If the input is not evenly divisible by three, it is right-padded by one or two octets of zeroes, and one or two "=" symbols are appended to the output.

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P |
| 1 | Q | R | S | T | U | V | W | X | Y | Z | a | b | c | d | e | f |
| 2 | g | h | i | j | k | l | m | n | o | p | q | r | s | t | u | v |
| 3 | w | x | y | z | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | + | / |
| The "=" sign is used for padding. | | | | | | | | | | | | | | | | |

- Base64 is position dependent.
- In general, any given string will be encoded in one of three different ways, depending on its offset into the input.
- This is what makes the encoding of patterns so hard.

| Encoding of "foo" at Three Different Offsets | |
|---|---|
| Offset 0 | Zm9v |
| Offset 1 | [159BFJNRVZdhlptx]mb2[+/8-9] |
| Offset 2 | [2GWm]Zvb[+/0-9w-z] |

DVLabs

- Our choice of pattern language controls the overall complexity of patterns for which we can search.
- If we were just looking for static strings, we wouldn't need a new tool - just use grep.
- Regular expressions provide a good balance of ease of implementation, expressive power, and common availability.

- Formalized by Stephen Kleene in 1956.
- Ken Thompson incorporated them as a useful pattern matching tool into his version of the QED editor for MIT's CTSS timesharing system.
- This later influenced Ken Thompson's implementation of ed for UNIX.
- From UNIX, regular expressions spread around the world.

- Character matches (e.g. "a")
- Concatenation (e.g. "ab")
- Alternation (e.g. "(ab|cd)")
- Character classes and inverse classes (e.g. "$[0-9]$")
- Kleene closures (e.g. "A*C")

- Backreferences and captures
- Variable-length repetition
- Left and right anchors
- Just about everything else

- Come up with a way to transform an arbitrary regular expression such that it will match its input when that input has been encoded using the Base64 algorithm.

- Do this transformation in such a way that the regular expression does not grow too large.

- Do this transformation in such a way that not too much information is lost.

- Do this transformation in such a way that the expression will match regardless of the pattern's offset from the beginning of input.

- Come up with a way to transform an arbitrary regular expression such that it will match its input when that input has been encoded using the Base64 algorithm.
- Do this transformation in such a way that the regular expression does not grow too large.
- Do this transformation in such a way that not too much information is lost.
- Do this transformation in such a way that the expression will match regardless of the pattern's offset from the beginning of input.

- Come up with a way to transform an arbitrary regular expression such that it will match its input when that input has been encoded using the Base64 algorithm.
- Do this transformation in such a way that the regular expression does not grow too large.
- Do this transformation in such a way that not too much information is lost.
- Do this transformation in such a way that the expression will match regardless of the pattern's offset from the beginning of input.

- Come up with a way to transform an arbitrary regular expression such that it will match its input when that input has been encoded using the Base64 algorithm.
- Do this transformation in such a way that the regular expression does not grow too large.
- Do this transformation in such a way that not too much information is lost.
- Do this transformation in such a way that the expression will match regardless of the pattern's offset from the beginning of input.

- Convert the regular expression into a nondeterministic finite state automaton using Thompson's algorithm, then convert the NFA to a deterministic finite state automaton using the powerset construction algorithm, then transform the DFA into a directed acyclic graph, then perform graph reductions until the graph is in a minimal form, transform the minimal form using graph transformations, then reduce again, then serialize the graph as a regular expression.
- This is hard, and I'm lazy.

- Convert the regular expression into a nondeterministic finite state automaton using Thompson's algorithm, then convert the NFA to a deterministic finite state automaton using the powerset construction algorithm, then transform the DFA into a directed acyclic graph, then perform graph reductions until the graph is in a minimal form, transform the minimal form using graph transformations, then reduce again, then serialize the graph as a regular expression.

- This is hard, and I'm lazy.

- Convert the regular expression into a nondeterministic finite state automaton using Thompson's algorithm, then convert the NFA to a deterministic finite state automaton using the powerset construction algorithm, then transform the DFA into a directed acyclic graph, then perform graph reductions until the graph is in a minimal form, transform the minimal form using graph transformations, then reduce again, then serialize the graph as a regular expression.

- This is hard, and I'm lazy.

- Enumerate all possible matching strings for a regular expression.

- Encode these strings and concatenate them inside an alternating regular expression.

- This would definitely work, with one problem...

- Regular expressions with Kleene closures can match an infinite number of strings.

- Enumerating an infinite number of strings can take a really long time.

- Even the trivial regular expression "...." would match over four billion strings.

- Enumerate all possible matching strings for a regular expression.

- Encode these strings and concatenate them inside an alternating regular expression.

- This would definitely work, with one problem...

- Regular expressions with Kleene closures can match an infinite number of strings.

- Enumerating an infinite number of strings can take a really long time.

- Even the trivial regular expression "...." would match over four billion strings.

- Enumerate all possible matching strings for a regular expression.
- Encode these strings and concatenate them inside an alternating regular expression.
- This would definitely work, with one problem...
- Regular expressions with Kleene closures can match an infinite number of strings.
- Enumerating an infinite number of strings can take a really long time.
- Even the trivial regular expression "...." would match over four billion strings.

- Enumerate all possible matching strings for a regular expression.
- Encode these strings and concatenate them inside an alternating regular expression.
- This would definitely work, with one problem...
- Regular expressions with Kleene closures can match an infinite number of strings.
- Enumerating an infinite number of strings can take a really long time.
- Even the trivial regular expression "...." would match over four billion strings.

**DVLabs**

- Enumerate all possible matching strings for a regular expression.
- Encode these strings and concatenate them inside an alternating regular expression.
- This would definitely work, with one problem...
- Regular expressions with Kleene closures can match an infinite number of strings.
- Enumerating an infinite number of strings can take a really long time.
- Even the trivial regular expression "...." would match over four billion strings.

- Enumerate all possible matching strings for a regular expression.
- Encode these strings and concatenate them inside an alternating regular expression.
- This would definitely work, with one problem...
- Regular expressions with Kleene closures can match an infinite number of strings.
- Enumerating an infinite number of strings can take a really long time.
- Even the trivial regular expression "...." would match over four billion strings.

- Enumerate all possible matching strings for a regular expression.
- Encode these strings and concatenate them inside an alternating regular expression.
- This would definitely work, with one problem...
- Regular expressions with Kleene closures can match an infinite number of strings.
- Enumerating an infinite number of strings can take a really long time.
- Even the trivial regular expression "...." would match over four billion strings.

- Enumerating all possible matching strings would, in theory, give us the correct answer.
- Instead of enumerating all possible strings and encoding them, what if we encoded the operations?
- In other words, what if we enumerated only what could match a given operation at a given point in time?

- Enumerating all possible matching strings would, in theory, give us the correct answer.

- Instead of enumerating all possible strings and encoding them, what if we encoded the operations?

- In other words, what if we enumerated only what could match a given operation at a given point in time?

- Enumerating all possible matching strings would, in theory, give us the correct answer.

- Instead of enumerating all possible strings and encoding them, what if we encoded the operations?

- In other words, what if we enumerated only what could match a given operation at a given point in time?

- Enumerating all possible matching strings would, in theory, give us the correct answer.
- Instead of enumerating all possible strings and encoding them, what if we encoded the operations?
- In other words, what if we enumerated only what could match a given operation at a given point in time?

**DVLabs**

# THE ALGORITHM
In a nutshell...

SCREAM

Rob King

Introduction

The Problem
The Base64
Algorithm
Regular Expressions

The Algorithm
Ways of Solving the
Problem
Encoding Operations

Performance

Expression
Optimization
Performance
Analysis

Implementation
and Usage
Common Use Cases
Caveats

Summary

- Create a list of regular expressions that match a fixed-length string.

- Enumerate the bitstrings matched by each of these expressions.

- Break each of these bitstrings into six-bit units.

- Encode each of these units.

- Treat each of these encoded strings as a branch in an $n$-way alternation in a regular expression.

- Optimize the expression.

- There is special handling for the two "meta" operations: alternations and closures.

- Create a list of regular expressions that match a fixed-length string.

- Enumerate the bitstrings matched by each of these expressions.

- Break each of these bitstrings into six-bit units.

- Encode each of these units.

- Treat each of these encoded strings as a branch in an $n$-way alternation in a regular expression.

- Optimize the expression.

- There is special handling for the two "meta" operations: alternations and closures.

- Create a list of regular expressions that match a fixed-length string.

- Enumerate the bitstrings matched by each of these expressions.

- Break each of these bitstrings into six-bit units.

- Encode each of these units.

- Treat each of these encoded strings as a branch in an *n*-way alternation in a regular expression.

- Optimize the expression.

- There is special handling for the two "meta" operations: alternations and closures.

- Create a list of regular expressions that match a fixed-length string.
- Enumerate the bitstrings matched by each of these expressions.
- Break each of these bitstrings into six-bit units.
- Encode each of these units.
- Treat each of these encoded strings as a branch in an $n$-way alternation in a regular expression.
- Optimize the expression.
- There is special handling for the two "meta" operations: alternations and closures.

- Create a list of regular expressions that match a fixed-length string.
- Enumerate the bitstrings matched by each of these expressions.
- Break each of these bitstrings into six-bit units.
- Encode each of these units.
- Treat each of these encoded strings as a branch in an $n$-way alternation in a regular expression.
- Optimize the expression.
- There is special handling for the two "meta" operations: alternations and closures.

- Create a list of regular expressions that match a fixed-length string.
- Enumerate the bitstrings matched by each of these expressions.
- Break each of these bitstrings into six-bit units.
- Encode each of these units.
- Treat each of these encoded strings as a branch in an *n*-way alternation in a regular expression.
- Optimize the expression.
- There is special handling for the two "meta" operations: alternations and closures.

- Create a list of regular expressions that match a fixed-length string.
- Enumerate the bitstrings matched by each of these expressions.
- Break each of these bitstrings into six-bit units.
- Encode each of these units.
- Treat each of these encoded strings as a branch in an *n*-way alternation in a regular expression.
- Optimize the expression.
- There is special handling for the two "meta" operations: alternations and closures.

- Create a list of regular expressions that match a fixed-length string.
- Enumerate the bitstrings matched by each of these expressions.
- Break each of these bitstrings into six-bit units.
- Encode each of these units.
- Treat each of these encoded strings as a branch in an *n*-way alternation in a regular expression.
- Optimize the expression.
- There is special handling for the two "meta" operations: alternations and closures.

- Non-fixed-length operations leave ambiguity as to what bits to place in a given six-bit unit.
- The goal of this phase of the algorithm is to get a collection of regular expressions that match fixed-length strings.
- First, let's go over how to encode the atomic operations.

This is a fairly obvious - a single character is encoded as the bitstring for that character.

Character classes are stored as a list of all characters
included in the class. This is not the most efficient possible
choice in terms of memory usage, but its implementation is
simpler and it is often faster.

| [ABC] | | | | | | | |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 |

Inverted character classes simply store the complement of
the list of indicated characters.

Wildcards are simply encoded as a character class with 256 entries. In the underlying implementation they are encoded as a special atom to save space, but from the point of view of the algorithm there is no difference.

- Single characters, character classes, and wildcards are treated as the building blocks of lists.
- Each fixed-length expression is simply a list of these operations.
- Therefore, concatenation is modelled implicitly by the ordering of these operations in the list.

- It was stated earlier that the algorithm deals only with regular expressions that match only fixed-length strings.
- Alternations can easily violate this: if one branch of an alternation has more character matches than the other, the expression overall is not fixed length.

- It was stated earlier that the algorithm deals only with regular expressions that match only fixed-length strings.
- Alternations can easily violate this: if one branch of an alternation has more character matches than the other, the expression overall is not fixed length.

- It was stated earlier that the algorithm deals only with regular expressions that match only fixed-length strings.
- Alternations can easily violate this: if one branch of an alternation has more character matches than the other, the expression overall is not fixed length.

- Alternations turn our simple lists of operations into directed acyclic graphs.
- For example, below is the representation of the regular expression "A(B|(C|D)E)F".



A(B|(C|D)E)F

- We can easily create a list of fixed-length expressions by enumerating all topological sorts of the graph.
- For each distinct topological sort, we can create a list of fixed-length operations.
- This reduces an alternation to a fixed-length structure, and can be performed recursively, resulting in a list of expressions, each of a fixed-length.

A(B|(C|D)E)F

- By adding an explicit "end of expression" pseudo-operation to the end of every expression, we can guarantee that all expressions have a shared tail.
- We can also add some number of explicit "start of expression" pseudo-operations from every entry point into the expression.
- These operations does not actually match anything; they exist solely for ease of implementation.



A(B|(C|D)E)F

- By finding a path from every START node to the END node, we can produce a list of fixed-length expressions.
- Note that this is essentially equivalent of the list of topological sorts, but provides the storage advantage of a shared tail.
- Implementation is also considerably easier and performs much better, since we don't need to implement a general topological sort; rather we can simply maintain a list of all start nodes and follow the only eminating edge from every node until we reach the end.

- We can apply this recursively, to build up a list of fixed-length expressions from a regular expression, even with nested alternations.
- Each of these fixed length expressions can be encoded easily using the methods for the fixed-length operations above.

- Using the algorithm above, we can create a list of fixed-length expressions consisting of eight bit units (for single characters) or lists of eight bit units (for character classes and wildcards).

- We now use a consumer function to walk through each fixed-length expression, extracting six bits at a time.
- This produces a new list of six-bit units.
- When the consumer function must get bits from a single character and a character class, or two classes, a list of six bit units is placed in the result list.

**DVLabs**

- Since the expression is going to be examining data in a streaming context, left and right anchors don't make sense in most situations.

- Therefore, if any bits are needed to complete a six-bit unit at the beginning or the end of the expression, they are treated as if a wildcard were present immediately before or after the expression.

- Since the expression is going to be examining data in a streaming context, left and right anchors don't make sense in most situations.
- Therefore, if any bits are needed to complete a six-bit unit at the beginning or the end of the expression, they are treated as if a wildcard were present immediately before or after the expression.

- Since the expression is going to be examining data in a streaming context, left and right anchors don't make sense in most situations.

- Therefore, if any bits are needed to complete a six-bit unit at the beginning or the end of the expression, they are treated as if a wildcard were present immediately before or after the expression.

**DVLabs**

- Each six bit unit is then replaced with its equivalent symbol in the Base64 alphabet.
- For lists of six bit units (from character classes), a list of Base64 symbols is produced.

- The list of fixed-length encoded expressions is joined into an *n*-way alternation, which is treated as a regular expression.
- This expression is first run through the optimizer, which recursively refactors common prefixes and suffixes.
- The expression is then reprocessed two more times, to account for varying offsets from the beginning of input.

**DVLabs**

SCREAM

Rob King

Introduction
The Problem
The Base64 Algorithm
Regular Expressions

The Algorithm
Ways of Solving the Problem
Encoding Operations

Performance
Expression Optimization
Performance Analysis

Implementation and Usage
Common Use Cases
Caveats

Summary

- Expressions with Kleene closures violate the fixed-length rule, since subexpressions can appear from zero to infinitely many times.

- Expressions with Kleene closures violate the acylic rule, since repeating expressions would become cycles in the graph representation.

**DVLabs**

- The easiest case for encoding closures is the case where they don't appear at all.

- Whenever a closure is present in an expression, it acts as a virtual alternation, with one branch having at least one instance of the closure, and the other branch having no instances.

- Closures are encoded three times, to account for varying offsets from the beginning of input.
- These three encoded expressions are placed in a special data structure that marks them as the components of a closure.

- Expressions with closures are treated as a list of expressions.
- The expressions are split such that each specially-encoded closure is a separate expression.
- Lists of expressions are encoded such that the implicit wildcard at the beginning and ending of each expression is replaced with the appropriate values from the next or previous expression in the list.

- The expression is treated as a list of constant-length expressions, from each possible starting expression to the logical end.

- All of these expressions will differ only around places where alternations were present.

- Therefore, everything leading up to the alternation, and everything after the alternation, will be identical for various expressions.

- We can therefore optimize the expression by refactoring all common prefixes and suffixes.

- The expression is treated as a list of constant-length expressions, from each possible starting expression to the logical end.
- All of these expressions will differ only around places where alternations were present.
- Therefore, everything leading up to the alternation, and everything after the alternation, will be identical for various expressions.
- We can therefore optimize the expression by refactoring all common prefixes and suffixes.

- The expression is treated as a list of constant-length expressions, from each possible starting expression to the logical end.

- All of these expressions will differ only around places where alternations were present.

- Therefore, everything leading up to the alternation, and everything after the alternation, will be identical for various expressions.

- We can therefore optimize the expression by refactoring all common prefixes and suffixes.

- The expression is treated as a list of constant-length expressions, from each possible starting expression to the logical end.

- All of these expressions will differ only around places where alternations were present.

- Therefore, everything leading up to the alternation, and everything after the alternation, will be identical for various expressions.

- We can therefore optimize the expression by refactoring all common prefixes and suffixes.

- The expression is treated as a list of constant-length expressions, from each possible starting expression to the logical end.
- All of these expressions will differ only around places where alternations were present.
- Therefore, everything leading up to the alternation, and everything after the alternation, will be identical for various expressions.
- We can therefore optimize the expression by refactoring all common prefixes and suffixes.

- In general, the algorithm will run in $O(n^2)$ time, since every addition of an alternation results in twice as many START nodes.

- This can of course become rapidly intractable, since even having 32 alternations would result in over four billion START nodes.

- While we do perform some optimizations to improve the theoretical running time, it is useful to point out that in the vast majority of practical runs, expressions generally have fewer than twenty alternations.

- In general, the algorithm will run in $O(n^2)$ time, since every addition of an alternation results in twice as many START nodes.

- This can of course become rapidly intractable, since even having 32 alternations would result in over four billion START nodes.

- While we do perform some optimizations to improve the theoretical running time, it is useful to point out that in the vast majority of practical runs, expressions generally have fewer than twenty alternations.

- In general, the algorithm will run in $O(n^2)$ time, since every addition of an alternation results in twice as many START nodes.

- This can of course become rapidly intractable, since even having 32 alternations would result in over four billion START nodes.

- While we do perform some optimizations to improve the theoretical running time, it is useful to point out that in the vast majority of practical runs, expressions generally have fewer than twenty alternations.

- In general, the algorithm will run in $O(n^2)$ time, since every addition of an alternation results in twice as many START nodes.
- This can of course become rapidly intractable, since even having 32 alternations would result in over four billion START nodes.
- While we do perform some optimizations to improve the theoretical running time, it is useful to point out that in the vast majority of practical runs, expressions generally have fewer than twenty alternations.

# PERFORMANCE ANALYSIS
Optimizing the Algorithm

SCREAM

Rob King

Introduction
The Problem
The Base64
Algorithm
Regular Expressions

The Algorithm
Ways of Solving the
Problem
Encoding Operations

Performance
Expression
Optimization
Performance
Analysis

Implementation
and Usage
Common Use Cases
Caveats

Summary

- The primary space optimization is the shared tail of the core data structure. Outside of pathological cases, this saves considerable storage.

- The alternation processing portion of the algorithm in reality only splits the expression when the two branches of the expression are of unequal length. This reduces the number of paths in many situations.

- Most of the calculations are memoized; that is, they are run only once for any given input for any given position and offset.

- With all of these optimizations taken into account, for most regular expressions in our test data set, runtimes of under three minutes are found (though memory usage is roughly quadratic on the number of operations in the expression).

- The primary space optimization is the shared tail of the core data structure. Outside of pathological cases, this saves considerable storage.

- The alternation processing portion of the algorithm in reality only splits the expression when the two branches of the expression are of unequal length. This reduces the number of paths in many situations.

- Most of the calculations are memoized; that is, they are run only once for any given input for any given position and offset.

- With all of these optimizations taken into account, for most regular expressions in our test data set, runtimes of under three minutes are found (though memory usage is roughly quadratic on the number of operations in the expression).

- The primary space optimization is the shared tail of the core data structure. Outside of pathological cases, this saves considerable storage.
- The alternation processing portion of the algorithm in reality only splits the expression when the two branches of the expression are of unequal length. This reduces the number of paths in many situations.
- Most of the calculations are memoized; that is, they are run only once for any given input for any given position and offset.
- With all of these optimizations taken into account, for most regular expressions in our test data set, runtimes of under three minutes are found (though memory usage is roughly quadratic on the number of operations in the expression).

- The primary space optimization is the shared tail of the core data structure. Outside of pathological cases, this saves considerable storage.
- The alternation processing portion of the algorithm in reality only splits the expression when the two branches of the expression are of unequal length. This reduces the number of paths in many situations.
- Most of the calculations are memoized; that is, they are run only once for any given input for any given position and offset.
- With all of these optimizations taken into account, for most regular expressions in our test data set, runtimes of under three minutes are found (though memory usage is roughly quadratic on the number of operations in the expression).

- The primary space optimization is the shared tail of the core data structure. Outside of pathological cases, this saves considerable storage.
- The alternation processing portion of the algorithm in reality only splits the expression when the two branches of the expression are of unequal length. This reduces the number of paths in many situations.
- Most of the calculations are memoized; that is, they are run only once for any given input for any given position and offset.
- With all of these optimizations taken into account, for most regular expressions in our test data set, runtimes of under three minutes are found (though memory usage is roughly quadratic on the number of operations in the expression).

- The initial version of the tool was written in Lua.

- However, the tool was quickly rewritten in Erlang.

- Erlang provided a much cleaner method of implementing shared-tail lists, memoization, and bit manipulation.

  - Total implementation is around 1500 lines of heavily-commented Erlang.

  - A large portion of this consists of a hand-crafted PEG library replacing the original YECC based parser.

- The initial version of the tool was written in Lua.

- However, the tool was quickly rewritten in Erlang.

- Erlang provided a much cleaner method of implementing shared-tail lists, memoization, and bit manipulation.

    - Total implementation is around 1500 lines of heavily-commented Erlang.

    - A large portion of this consists of a hand-crafted PEG library replacing the original YECC based parser.

- The initial version of the tool was written in Lua.

- However, the tool was quickly rewritten in Erlang.

- Erlang provided a much cleaner method of implementing shared-tail lists, memoization, and bit manipulation.
  - Total implementation is around 1500 lines of heavily-commented Erlang.
  - A large portion of this consists of a hand-crafted PEG library replacing the original YECC based parser.

- The initial version of the tool was written in Lua.
- However, the tool was quickly rewritten in Erlang.
- Erlang provided a much cleaner method of implementing shared-tail lists, memoization, and bit manipulation.
  - Total implementation is around 1500 lines of heavily-commented Erlang.
  - A large portion of this consists of a hand-crafted PEG library replacing the original YECC based parser.

- Erlang makes it very easy to parallelize computation across multiple threads, processors, and machines.
- The algorithm itself is very amenable to parallelization, and an experimental version of the tool has shown significant speedups.
- This experimental version was also compiled using Erlang's "HiPE" native-code compiler where possible. This resulted in an additional 20% on average speed increase.

- Erlang makes it very easy to parallelize computation across multiple threads, processors, and machines.
- The algorithm itself is very amenable to parallelization, and an experimental version of the tool has shown significant speedups.
- This experimental version was also compiled using Erlang's "HiPE" native-code compiler where possible. This resulted in an additional 20% on average speed increase.

**DVLabs**

- Erlang makes it very easy to parallelize computation across multiple threads, processors, and machines.
- The algorithm itself is very amenable to parallelization, and an experimental version of the tool has shown significant speedups.
- This experimental version was also compiled using Erlang's "HiPE" native-code compiler where possible. This resulted in an additional 20% on average speed increase.

- While a traditional store-and-forward email server is of course optimal for filtering email, b64re-based signatures on an IPS have been useful in prefiltering email.

- One common use case is filtering obviously malicious email, based on shellcode detection or other parameters, lowering the load on the backend mail server.

- Another interesting application is protecting store-and-forward scanners from themselves. Some systems have bugs that can be triggered by specially formatted emails; an inline IPS can filter these attacks.

- While a traditional store-and-forward email server is of course optimal for filtering email, b64re-based signatures on an IPS have been useful in prefiltering email.

- One common use case is filtering obviously malicious email, based on shellcode detection or other parameters, lowering the load on the backend mail server.

- Another interesting application is protecting store-and-forward scanners from themselves. Some systems have bugs that can be triggered by specially formatted emails; an inline IPS can filter these attacks.

SCREAM

Rob King

Introduction
The Problem
The Base64
Algorithm
Regular Expressions

The Algorithm
Ways of Solving the
Problem
Encoding Operations

Performance
Expression
Optimization
Performance
Analysis

Implementation
and Usage
Common Use Cases
Caveats

Summary

- While a traditional store-and-forward email server is of course optimal for filtering email, b64re-based signatures on an IPS have been useful in prefiltering email.

- One common use case is filtering obviously malicious email, based on shellcode detection or other parameters, lowering the load on the backend mail server.

- Another interesting application is protecting store-and-forward scanners from themselves. Some systems have bugs that can be triggered by specially formatted emails; an inline IPS can filter these attacks.

- While a traditional store-and-forward email server is of course optimal for filtering email, b64re-based signatures on an IPS have been useful in prefiltering email.

- One common use case is filtering obviously malicious email, based on shellcode detection or other parameters, lowering the load on the backend mail server.

- Another interesting application is protecting store-and-forward scanners from themselves. Some systems have bugs that can be triggered by specially formatted emails; an inline IPS can filter these attacks.

- Email is a common source of sensitive data leakage.
- An inline IPS running signatures using this tool can be used to scan for outbound sensitive information, even if an external mail server is used.

- Email is a common source of sensitive data leakage.

- An inline IPS running signatures using this tool can be used to scan for outbound sensitive information, even if an external mail server is used.

- Email is a common source of sensitive data leakage.
- An inline IPS running signatures using this tool can be used to scan for outbound sensitive information, even if an external mail server is used.

- Expressions processed in this way are inherently looser than their unencoded counterparts. Why?

- Recall that the expression is encoded three times, to account for varying offsets from the beginning of input.

- The expression designed to match when its offset is $n + 1$ from the beginning of the input could match input that is actually at a different modular offset.

- This is mitigated somewhat if the expression can be anchored by long constant strings, or if the beginning of input is known.

DVLabs

- Expressions processed in this way are inherently looser than their unencoded counterparts. Why?

- Recall that the expression is encoded three times, to account for varying offsets from the beginning of input.

- The expression designed to match when its offset is $n + 1$ from the beginning of the input could match input that is actually at a different modular offset.

- This is mitigated somewhat if the expression can be anchored by long constant strings, or if the beginning of input is known.

- Expressions processed in this way are inherently looser than their unencoded counterparts. Why?
- Recall that the expression is encoded three times, to account for varying offsets from the beginning of input.
- The expression designed to match when its offset is $n + 1$ from the beginning of the input could match input that is actually at a different modular offset.
- This is mitigated somewhat if the expression can be anchored by long constant strings, or if the beginning of input is known.

SCREAM

Rob King

Introduction
The Problem
The Base64 Algorithm
Regular Expressions

The Algorithm
Ways of Solving the Problem
Encoding Operations

Performance
Expression Optimization
Performance Analysis

Implementation and Usage
Common Use Cases
Caveats

Summary

- Expressions processed in this way are inherently looser than their unencoded counterparts. Why?
- Recall that the expression is encoded three times, to account for varying offsets from the beginning of input.
- The expression designed to match when its offset is $n + 1$ from the beginning of the input could match input that is actually at a different modular offset.
- This is mitigated somewhat if the expression can be anchored by long constant strings, or if the beginning of input is known.

- Expressions processed in this way are inherently looser than their unencoded counterparts. Why?
- Recall that the expression is encoded three times, to account for varying offsets from the beginning of input.
- The expression designed to match when its offset is $n + 1$ from the beginning of the input could match input that is actually at a different modular offset.
- This is mitigated somewhat if the expression can be anchored by long constant strings, or if the beginning of input is known.

**DVLabs**

- Most Base64 processors allow the insertion of newlines at arbitrary points of the input.
- The encoded expression can be modified to include an optional newline after every character, though this is expensive.
- Some regular expression engines offer the ability to strip newlines from input, removing the burden on the expression.

- Most Base64 processors allow the insertion of newlines at arbitrary points of the input.

- The encoded expression can be modified to include an optional newline after every character, though this is expensive.

- Some regular expression engines offer the ability to strip newlines from input, removing the burden on the expression.

- Most Base64 processors allow the insertion of newlines at arbitrary points of the input.

- The encoded expression can be modified to include an optional newline after every character, though this is expensive.

- Some regular expression engines offer the ability to strip newlines from input, removing the burden on the expression.

- Most Base64 processors allow the insertion of newlines at arbitrary points of the input.

- The encoded expression can be modified to include an optional newline after every character, though this is expensive.

- Some regular expression engines offer the ability to strip newlines from input, removing the burden on the expression.

**DVLabs**

- Using this tool, users can reap several benefits:
  - Pre-filter encoded data streams, lessening the need for expensive decoding.
  - Protect infrastructure that provides decoding services but is itself vulnerable
  - Inspect encoded streams within incurring the overhead of decoding.
- This talk illustrates some interesting challenges when inspecting encoded streams of data, especially in a performance-critical way, and gives some possible solutions.

- Using this tool, users can reap several benefits:
  - Pre-filter encoded data streams, lessening the need for expensive decoding.
  - Protect infrastructure that provides decoding services but is itself vulnerable
  - Inspect encoded streams within incurring the overhead of decoding.
- This talk illustrates some interesting challenges when inspecting encoded streams of data, especially in a performance-critical way, and gives some possible solutions.

- Using this tool, users can reap several benefits:
  - Pre-filter encoded data streams, lessening the need for expensive decoding.
  - Protect infrastructure that provides decoding services but is itself vulnerable
  - Inspect encoded streams within incurring the overhead of decoding.
- This talk illustrates some interesting challenges when inspecting encoded streams of data, especially in a performance-critical way, and gives some possible solutions.

**DVLabs**

- Rob at Work: rking@tippingpoint.com
- Rob at Home: jking@deadpixi.com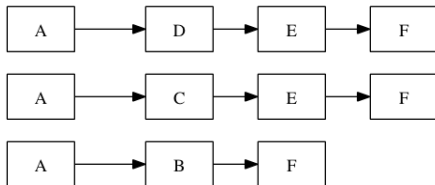