

*Arjan van Leeuwen*

 *@avl7771*



# Dealing with Strings in C++





# The Story of RSS and the StringHash



**“When I check for new  
feed items, Opera  
*almost* crashes”**



# An RSS feed

```
<?xml version="1.0" encoding="UTF-8"?>
<rss version="2.0">
  <channel>
    <title>Example RSS Feed</title>
    <item>
      <title>Example Item</title>
      <description>A summary.</description>
      <link>http://www.example.com/foo</link>
      <pubDate>Mon, 23 Sep 2013 03:00:05 GMT
      </pubDate>
    </item>
  </channel>
</rss>
```



# How long can it get?

Strings can sometimes get much longer than you expect, especially if it's some kind of user input being processed (which strings often are). We tend to forget strings have a variable length, and that the length of the string has an important effect on any operation done on it. Longer strings lead to all kinds of problems: slowness, unexpected results when



# Many operations $O(N)$

Operations	
<code>clear</code>	clears the contents (public member function)
<code>insert</code>	inserts characters (public member function)
<code>erase</code>	removes characters (public member function)
<code>push_back</code>	appends a character to the end (public member function)
<code>pop_back(C++11)</code>	removes the last character (public member function)
<code>append</code>	appends characters to the end (public member function)
<code>operator+=</code>	appends characters to the end (public member function)
<code>compare</code>	compares two strings (public member function)
<code>replace</code>	replaces specified portion of a string (public member function)
<code>substr</code>	returns a substring (public member function)
<code>copy</code>	copies characters (public member function)
<code>resize</code>	changes the number of characters stored (public member function)
<code>swap</code>	swaps the contents (public member function)
Search	
<code>find</code>	find characters in the string (public member function)
<code>rfind</code>	find the last occurrence of a substring (public member function)
<code>find_first_of</code>	find first occurrence of characters (public member function)
<code>find_first_not_of</code>	find first absence of characters (public member function)
<code>find_last_of</code>	find last occurrence of characters (public member function)
<code>find_last_not_of</code>	find last absence of characters (public member function)



**Do you actually need  
to use a string?**



# Mutable vs. immutable strings

A decorative graphic on the right side of the slide, consisting of a red mesh or grid pattern that curves and flows downwards, resembling a stylized ribbon or a piece of fabric.



# Mutable / Immutable

Language	String type	Mutable	Notes
<b>C</b>	None / char[]	Mutable	Change in length hard
<b>C++</b>	std::string	Mutable	
<b>Objective-C</b>	NSString	Immutable	NSMutableString available
<b>Python</b>	String	Immutable	
<b>JavaScript</b>	String	Immutable	



# Mutable vs. immutable strings

- Changes that change length are expensive
- Many strings never need to be changed
- Immutable views of strings and substrings can replace mutable strings in many common cases
  - e.g. tokenization, splitting, trimming, ...



# Problems with char arrays

- Is that a string or an array of 8-bit integers?
- `strlen()` is  $O(N)$
- NULL-terminated
- Unsafe (as any C array)



# The case for `std::string`

- No memory ownership issues
- Clear use case
- String length known
- Safe string utility functions
- Operators increase readability



YOU WOULDN'T COPY  
AN ARRAY

ASSIGNING STRINGS  
IS COPYING

COPYING  
IS AGAINST THE LAW!

# Problems with `std::string`

- Always allocates memory
- Makes copying invisible with assignment and copy-constructor
- String literals are always copied
- Any substring creates another copy
- Mutable type





**Does `std::string`  
make people lazy?**

*Or: Don't tell anybody, but I  
actually liked `char arrays`*



# Enter `std::string_view`

- ‘view on a string’: a non-owning string type
- String length known, safe helper functions
- Can be used with `std::string` and char arrays
- Copying to `std::string` is possible but explicit
- Real string literals!
- View can be a substring



# String types

	Type	Alloc	Copy	View
No STL	<code>char []</code>	Heap, stack	<code>strdup,</code> <code>str{n/l}cpy</code>	<code>const char*</code>
C++03	<code>std::string</code>	Heap	<code>std::string</code>	<code>const</code> <code>std::string&amp;</code>
C++17	<code>std::string_</code> <code>view</code>	Non- owning	<code>std::string</code>	<code>std::string_</code> <code>view</code>



# `std::string_view` everywhere!

- use a `std::string_view` wherever you can!
- When you're working with `std::strings`, take care of ownership to prevent copy
  - e.g. use `std::move` when content not needed



# Strings at compile-time



# Growing strings

- Be careful with modifying operations such as `append()`
- Avoid creating a string out of many parts, better to create at once (remember immutability)



# Growing strings

```
std::string CopyString(  
    const char* to_copy, size_t length) {  
    std::string copied;  
  
    for (size_t i = 0; i < length; i += BLOCKSIZE)  
        copied.append(to_copy + i,  
            std::min(BLOCKSIZE, length - i));  
  
    return copied;  
}
```



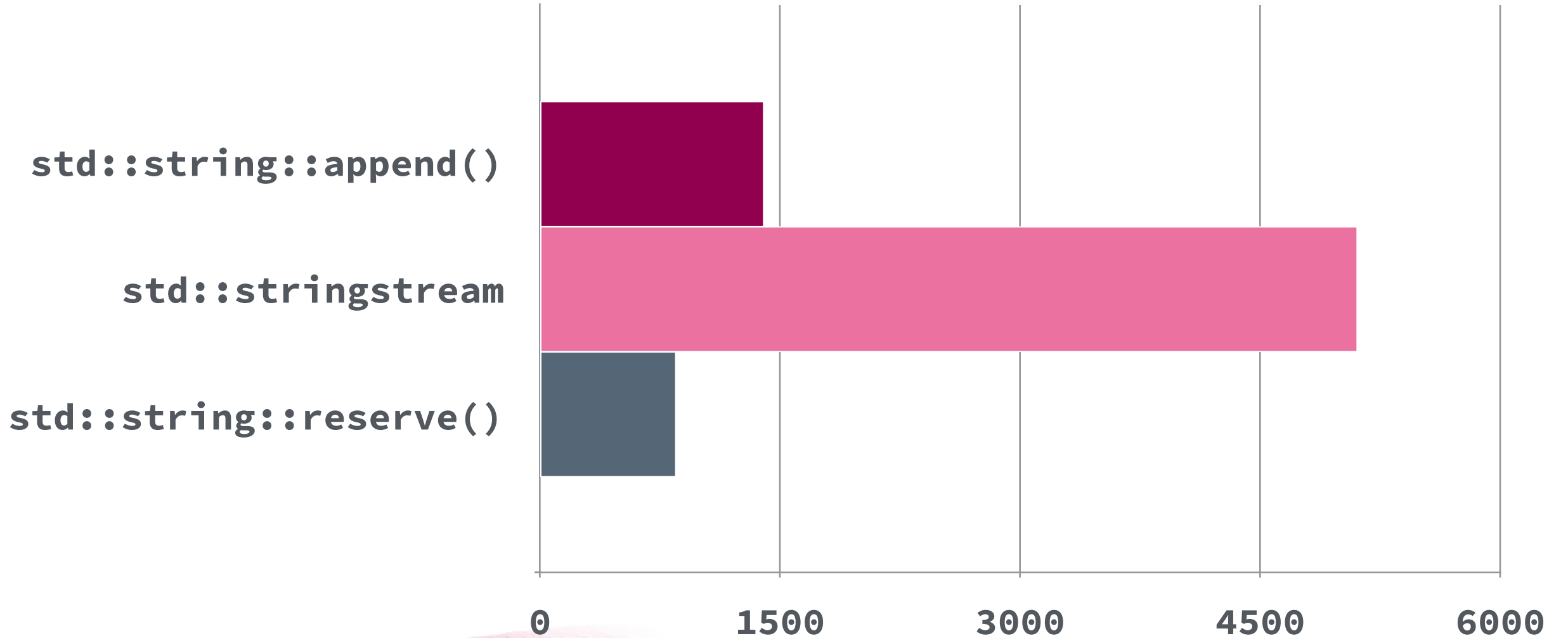
# Growing strings

```
std::string CopyString(  
    const char* to_copy, size_t length) {  
    std::stringstream copied;  
  
    for (size_t i = 0; i < length; i += BLOCKSIZE)  
        copied.write(to_copy + i,  
                    std::min(BLOCKSIZE, length - i));  
  
    return copied.str();  
}
```





# Growing strings



# Converting numbers

*When you need to switch  
between two worlds*



# Converting to string

```
std::string Convert(int i) {  
    std::stringstream stream;  
    stream << i;  
    return stream.str();  
}
```



# Converting to string

```
std::string Convert(int i) {  
    return std::to_string(i);  
}
```

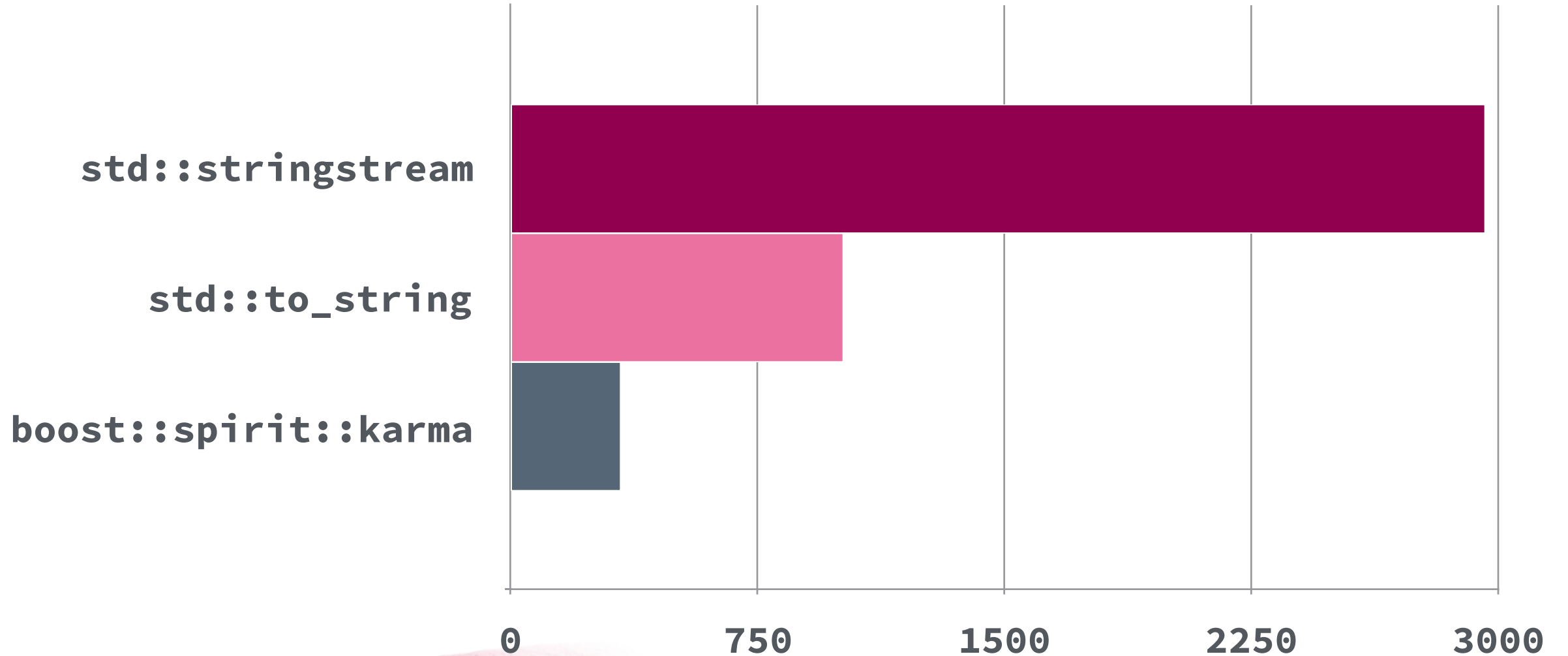


# Converting to string

```
std::string Convert(int i) {  
    namespace karma = boost::spirit::karma;  
    std::string converted;  
    std::back_inserter<std::string>  
        sink(converted);  
  
    karma::generate(sink, karma::int_, i);  
    return converted;  
}
```



# Converting to string



# Converting to integer

```
int Convert(const std::string& str) {  
    return std::stoi(str);  
}
```



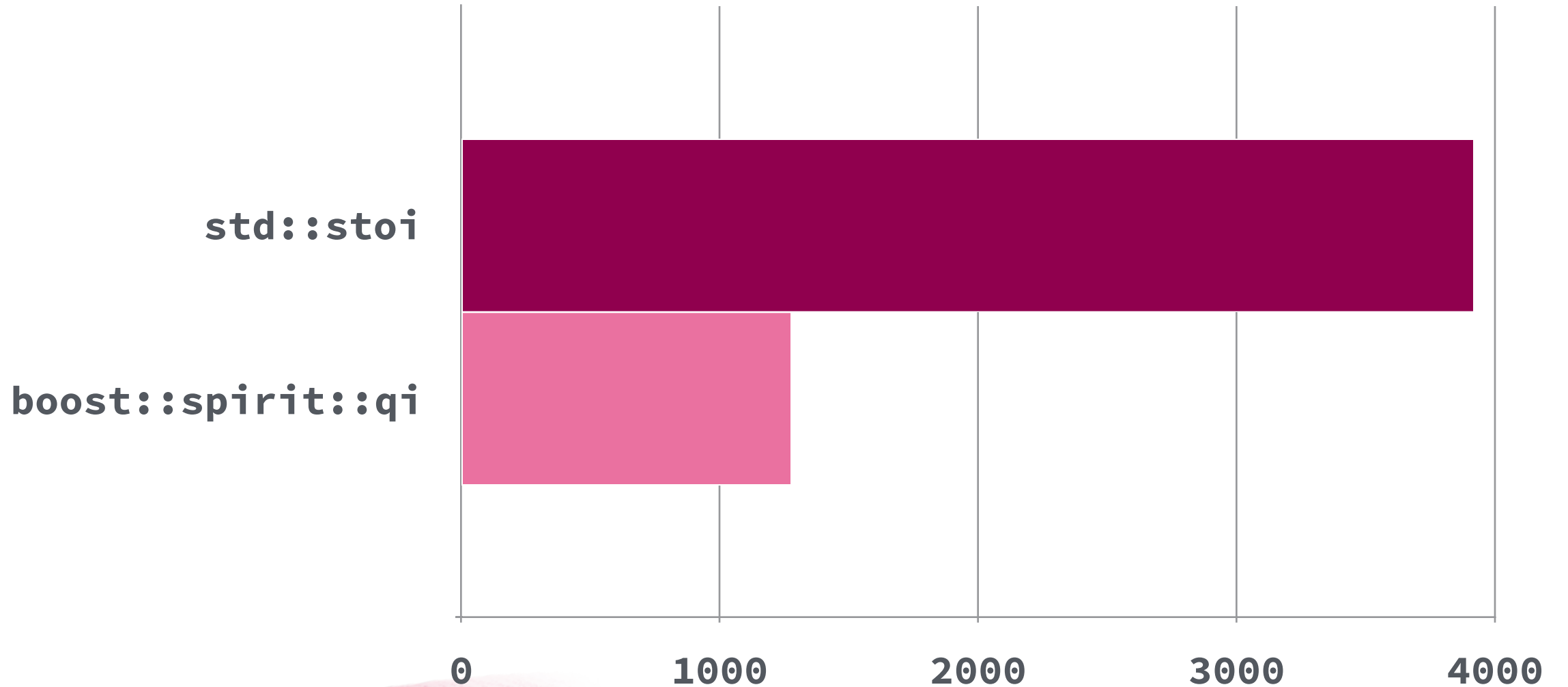
# Converting to integer

```
int Convert(const std::string& str) {  
    namespace qi = boost::spirit::qi;  
    int converted;  
  
    qi::parse(str.begin(), str.end(),  
              qi::int_, converted);  
    return converted;  
}
```





# Converting to integer



# Other options

`std::from_chars` *and* `std::to_chars`  
*are coming!*

low-level, locale-independent functions for conversions between integers and strings and between floating-point numbers and strings.



# Regular Expressions

*why the hotel  
really didn't have a pool*



# Avoid Regular Expressions

*if you can*



# Character encoding



# Basics of character encoding

- Language uses words or sentences made of characters
  - Latin 'á', Chinese '請' are examples of characters
- Code points are numbers assigned to characters
- Encoding specifies how to represent code points in bytes



# The Unicode standard

- Unicode specifies standard code points for many characters
  - Regularly updated to add emoji
- Unicode specifies several encodings that map all code points to byte sequences, called the Unicode Transformation Format (UTF)



# The Unicode standard

- Unicode specifies standard code points for many characters
  - Regularly updated to add emoji 🍌
- Unicode specifies several encodings that map all code points to byte sequences, called the Unicode Transformation Format (UTF)





# UTF encodings

	Size	Width per code point	Remarks
<b>UTF-8</b>	8 bit	Variable	Backwards compatible with ASCII
<b>UTF-16</b>	16 bit	Variable	Used in Windows system libraries, endianness matters
<b>UTF-32</b>	32 bit	Fixed	Used on Linux/BSD system libraries, endianness matters



# Why you should use UTF-8

- Backwards compatibility with ASCII is useful (no need to write separate functions)
- Uses less memory for ASCII-based data formats
- Uses more memory per code point for characters that actually represent more content
- No endianness concerns means compatibility



# When you don't use UTF-8

- Interacting continuously with APIs that require other encodings (e.g. UTF-16 in the Windows API)
- Can still use UTF-8 transcoding when saving files or doing network traffic
- Avoid transcoding as much as possible



# UTF encodings in C++

	Primitive type	String type	String view
UTF-8	<code>char</code>	<code>std::string</code>	<code>std::string_view</code>
UTF-16	<code>char16_t</code>	<code>std::u16string</code>	<code>std::u16string_view</code>
UTF-32	<code>char32_t</code>	<code>std::u32string</code>	<code>std::u32string_view</code>



# UTF encodings in C++

- Transcoding between different UTF variants is possible with the standard library but messy (look up `<codecvt>` for details)
- String literals can specify encoding and use non-ascii characters using code point specifiers:
  - `u8`"A 4 digit code point in UTF-8: **\u2018.**"
  - `u`"An 8 digit code point in UTF-16: **\U00002018.**"
  - `U`"An 8 digit code point in UTF-32: **\U00002018.**"



# **Dictionaries and indexes**



*Arjan van Leeuwen*

 *@avl7771*



# Dealing with Strings in C++

