# Data Compression and IoT

Ulises Tirado Zatarain [1]
(ulises.tirado@cimat.mx)

[1] Algorists Group

July, 2016

# Outline

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Outline

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

**Definitions**
Example
Some ideas and approaches

# What is data compression?

## Definition (Data compression)

Representation of information using less space than original data. The action to compress data is called **compression** and the opposite actions is called **decompression**. It's a particular case of encoding/decoding information.

- Kinds of compression:
  - Loseless
  - Lossy

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

## Loseless compression

- Information can be retrieved exactly as original data.
- Usually used for text compression
- Some known formats:
  - Zip
  - GZip
  - RAR
  - ACE
  - 7Zip
  - B2Zip
  - ...

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

**Definitions**
Example
Some ideas and approaches

# Lossy compression

- Information loses some data, that cannot be retrieved exactly as before it is compressed.

- Usually used for media compression: images, audio, video.

- Some known formats:
  - JPEG, GIF, PNG, ...
  - MP3, OGG, AAC, ...
  - H264, MPEG-4, VP8, ...

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

**Definitions**
Example
Some ideas and approaches

## Kinds of information

- **Redundant**
  - Repetitive data
  - Predictable data

- Irrelevant
  - Invisible data
  - Removing this data don't affect message content

- Basic
  - Essential data
  - It's needed to retrieve original data
  - It should be transmitted

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

**Definitions**
Example
Some ideas and approaches

# Kinds of information

- **Redundant**
  - Repetitive data
  - Predictable data

- Irrelevant
  - Invisible data
  - Removing this data don't affect message content

- Basic
  - Essential data
  - It's needed to retrieve original data
  - It should be transmitted

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

## Kinds of information

- **Redundant**
  - Repetitive data
  - Predictable data

- Irrelevant
  - Invisible data
  - Removing this data don't affect message content

- Basic
  - Essential data
  - It's needed to retrieve original data
  - It should be transmitted

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

**Definitions**
Example
Some ideas and approaches

## Kinds of information

- **Redundant**
  - Repetitive data
  - Predictable data

- **Irrelevant**
  - Invisible data
  - Removing this data don't affect message content

- Basic
  - Essential data
  - It's needed to retrieve original data
  - It should be transmitted

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Kinds of information

- **Redundant**
  - Repetitive data
  - Predictable data

- **Irrelevant**
  - Invisible data
  - Removing this data don't affect message content

- Basic
  - Essential data
  - It's needed to retrieve original data
  - It should be transmitted

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Kinds of information

- **Redundant**
  - Repetitive data
  - Predictable data

- **Irrelevant**
  - Invisible data
  - Removing this data don't affect message content

- Basic
  - Essential data
  - It's needed to retrieve original data
  - It should be transmitted

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Kinds of information

- **Redundant**
  - Repetitive data
  - Predictable data

- **Irrelevant**
  - Invisible data
  - Removing this data don't affect message content

- **Basic**
  - Essential data
  - It's needed to retrieve original data
  - It should be transmitted

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Kinds of information

- **Redundant**
  - Repetitive data
  - Predictable data

- **Irrelevant**
  - Invisible data
  - Removing this data don't affect message content

- **Basic**
  - Essential data
  - It's needed to retrieve original data
  - It should be transmitted

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

**Definitions**
Example
Some ideas and approaches

# Kinds of information

- **Redundant**
  - Repetitive data
  - Predictable data

- **Irrelevant**
  - Invisible data
  - Removing this data don't affect message content

- **Basic**
  - Essential data
  - It's needed to retrieve original data
  - It should be transmitted

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Kinds of information

- **Redundant**
  - Repetitive data
  - Predictable data

- **Irrelevant**
  - Invisible data
  - Removing this data don't affect message content

- **Basic**
  - Essential data
  - It's needed to retrieve original data
  - It should be transmitted

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Outline

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

## Lets see an example: Fruit 100% random & $\mathfrak{I}(\mathbf{country})$

There are six popular fruits in an imaginary random country with some states (about 32). People in the country implements an elections system to know: What's the favorite fruit ever in this random, imaginary and 100% hypothetical country?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

## Lets see an example: Fruit 100% random & $\mathfrak{I}(\textbf{country})$

There are six popular fruits in an imaginary random country with some states (about 32). People in the country implements an elections system to know: What's the favorite fruit ever in this random, imaginary and 100% hypothetical country?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Lets see an example: Fruit 100% random & $\Im(\textbf{country})$

There are six popular fruits in an imaginary random country with some states (about 32). People in the country implements an elections system to know: What's the favorite fruit ever in this random, imaginary and 100% hypothetical country?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Lets see an example: Fruit 100% random & $\mathfrak{I}(\mathbf{country})$

There are six popular fruits in an imaginary random country with some states (about 32). People in the country implements an elections system to know: What's the favorite fruit ever in this random, imaginary and 100% hypothetical country?



Can you see the different kinds of information?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

## Fruit 100% random & $\mathfrak{I}(\textbf{country})$: Game rules

The election system has following rules:

- Each citizen has an unique ID scanned from his/her ID card.

- Each citizen can vote only once and only by one fruit.

- If somebody tries to vote twice or more, then all votes from this citizen will be invalidated.

- Any citizen can vote in any state.

- There is a central system publishing partial live results.

- Each state has a system to votes counting and this reports to the central system. This systems only can report (to central system) votes from citizens who are natives from that state.

- In anytime the systems in each states can communicate with the other state systems to report votes from non-native citizens.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Fruit 100% random & $\mathfrak{I}(\mathbf{country})$: Game rules

The election system has following rules:

- Each citizen has an unique ID scanned from his/her ID card.
- Each citizen can vote only once and only by one fruit.
- If somebody tries to vote twice or more, then all votes from this citizen will be invalidated.
- Any citizen can vote in any state.
- There is a central system publishing partial live results.
- Each state has a system to votes counting and this reports to the central system. This systems only can report (to central system) votes from citizens who are natives from that state.
- In anytime the systems in each states can communicate with the other state systems to report votes from non-native citizens.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Fruit 100% random & $\mathfrak{I}(\mathbf{country})$: Game rules

The election system has following rules:

- Each citizen has an unique ID scanned from his/her ID card.
- Each citizen can vote only once and only by one fruit.
- If somebody tries to vote twice or more, then all votes from this citizen will be invalidated.
- Any citizen can vote in any state.
- There is a central system publishing partial live results.
- Each state has a system to votes counting and this reports to the central system. This systems only can report (to central system) votes from citizens who are natives from that state.
- In anytime the systems in each states can communicate with the other state systems to report votes from non-native citizens.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Fruit 100% random & $\mathfrak{I}(\mathbf{country})$: Game rules

The election system has following rules:

- Each citizen has an unique ID scanned from his/her ID card.
- Each citizen can vote only once and only by one fruit.
- If somebody tries to vote twice or more, then all votes from this citizen will be invalidated.
- Any citizen can vote in any state.
- There is a central system publishing partial live results.
- Each state has a system to votes counting and this reports to the central system. This systems only can report (to central system) votes from citizens who are natives from that state.
- In anytime the systems in each states can communicate with the other state systems to report votes from non-native citizens.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Fruit 100% random & $\mathfrak{I}(\mathbf{country})$: Game rules

The election system has following rules:

- Each citizen has an unique ID scanned from his/her ID card.
- Each citizen can vote only once and only by one fruit.
- If somebody tries to vote twice or more, then all votes from this citizen will be invalidated.
- Any citizen can vote in any state.
- There is a central system publishing partial live results.
- Each state has a system to votes counting and this reports to the central system. This systems only can report (to central system) votes from citizens who are natives from that state.
- In anytime the systems in each states can communicate with the other state systems to report votes from non-native citizens.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Fruit 100% random & $\mathfrak{I}(\textbf{country})$: Game rules

The election system has following rules:

- Each citizen has an unique ID scanned from his/her ID card.
- Each citizen can vote only once and only by one fruit.
- If somebody tries to vote twice or more, then all votes from this citizen will be invalidated.
- Any citizen can vote in any state.
- There is a central system publishing partial live results.
- Each state has a system to votes counting and this reports to the central system. This systems only can report (to central system) votes from citizens who are natives from that state.
- In anytime the systems in each states can communicate with the other state systems to report votes from non-native citizens.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Fruit 100% random & $\mathfrak{I}(\mathbf{country})$: Game rules

The election system has following rules:

- Each citizen has an unique ID scanned from his/her ID card.
- Each citizen can vote only once and only by one fruit.
- If somebody tries to vote twice or more, then all votes from this citizen will be invalidated.
- Any citizen can vote in any state.
- There is a central system publishing partial live results.
- Each state has a system to votes counting and this reports to the central system. This systems only can report (to central system) votes from citizens who are natives from that state.
- In anytime the systems in each states can communicate with the other state systems to report votes from non-native citizens.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Outline

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

## Logistics (brainstorming)

- How people can vote?
- Is an app needed?
- How people can vote outside of their state? (non-native people)
- Infrastructure?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

## Logistics (brainstorming)

- How people can vote?
- Is an app needed?
- How people can vote outside of their state? (non-native people)
- Infrastructure?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

## Logistics (brainstorming)

- How people can vote?
- Is an app needed?
- How people can vote outside of their state? (non-native people)
- Infrastructure?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Logistics (brainstorming)

- How people can vote?
- Is an app needed?
- How people can vote outside of their state? (non-native people)
- Infrastructure?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

## Architecture and design (brainstorming)

- First, think in the small case (i.e. one server by state)
- Solve for this case
- Improve to solve big case (i.e. dividing each states by districts)

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

## Architecture and design (brainstorming)

- First, think in the small case (i.e. one server by state)
- Solve for this case
- Improve to solve big case (i.e. dividing each states by districts)

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Architecture and design (brainstorming)

- First, think in the small case (i.e. one server by state)
- Solve for this case
- Improve to solve big case (i.e. dividing each states by districts)

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# What about data transferring? (brainstorming)

Do you have some ideas for the system?

- One vote once?
- Several votes at once?
- What technology can we use?
  - XML
  - JSON
  - Our own coding method?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# What about data transferring? (brainstorming)

Do you have some ideas for the system?

- One vote once?

- Several votes at once?

- What technology can we use?

  - XML
  - JSON
  - Our own coding method?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# What about data transferring? (brainstorming)

Do you have some ideas for the system?

- One vote once?

- Several votes at once?

- What technology can we use?

  - XML
  - JSON
  - Our own coding method?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# What about data transferring? (brainstorming)

Do you have some ideas for the system?

- One vote once?
- Several votes at once?
- What technology can we use?
    - XML
    - JSON
    - Our own coding method?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# What about data transferring? (brainstorming)

Do you have some ideas for the system?

- One vote once?

- Several votes at once?

- What technology can we use?

    - XML
    - JSON
    - Our own coding method?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# What about data transferring? (brainstorming)

Do you have some ideas for the system?

- One vote once?

- Several votes at once?

- What technology can we use?

    - XML
    - JSON
    - Our own coding method?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Data transferring: XML? (brainstorming)

```
1  <?xml version="1.0" encoding="UTF-8" ?>
2  <!DOCTYPE FruitCountry SYSTEM "votes.dtd">
3  <state id="25">
4      <vote>
5          <citizen id="111999" />
6          <by>Apple</by>
7      </vote>
8      ...
9      <vote>
10         <citizen id="333777" />
11         <by>Strawberry</by>
12     </vote>
13 </state>
```

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Data transferring: JSON? (brainstorming)

```
1  {
2      state: 25,
3      votes: [
4          { citizen: 111999, by: 'Apple' },
5          { citizen: 222888, by: 'Pear' },
6          { citizen: 222888, by: 'Banana' },
7          { citizen: 222888, by: 'Watermelon' },
8          ...
9          { citizen: 333777, by: 'Strawberry' }
10         { citizen: 333777, by: 'Orange' }
11     ]
12 }
```

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Data transferring: Our own coding method? (brainstorming)

- What if we use some abreviations?
  - A: Apple
  - B: Banana
  - O: Orange
  - P: Pear
  - S: Strawberry
  - W: Watermelon

- Do we really need to send the citizen ID?

- Do we really need to send the state ID?

- Fixed width messages?

- A possible message from state to central system:

$$25 \quad AAAAPPPPPBBBBBWWWSSOOOOOAAA$$

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Data transferring: Our own coding method? (brainstorming)

- What if we use some abreviations?
  - A: Apple
  - B: Banana
  - O: Orange
  - P: Pear
  - S: Strawberry
  - W: Watermelon

- Do we really need to send the citizen ID?

- Do we really need to send the state ID?

- Fixed width messages?

- A possible message from state to central system:

$$25 \quad AAAAPPPPPBBBBBWWWSSOOOOOAAA$$

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Data transferring: Our own coding method? (brainstorming)

- What if we use some abreviations?
  - A: Apple
  - B: Banana
  - O: Orange
  - P: Pear
  - S: Strawberry
  - W: Watermelon

- Do we really need to send the citizen ID?

- Do we really need to send the state ID?

- Fixed width messages?

- A possible message from state to central system:

$$25 \quad AAAAPPPPPBBBBBWWWSSOOOOOAAA$$

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Definitions
Example
Some ideas and approaches

# Data transferring: Our own coding method? (brainstorming)

- What if we use some abreviations?
  - A: Apple
  - B: Banana
  - O: Orange
  - P: Pear
  - S: Strawberry
  - W: Watermelon

- Do we really need to send the citizen ID?

- Do we really need to send the state ID?

- Fixed width messages?

- A possible message from state to central system:

    25   AAAAPPPPPBBBBBWWWSSOOOOOAAA

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

**Definitions**
Example
**Some ideas and approaches**

# Data transferring: Our own coding method? (brainstorming)

- What if we use some abreviations?
  - A: Apple
  - B: Banana
  - O: Orange
  - P: Pear
  - S: Strawberry
  - W: Watermelon

- Do we really need to send the citizen ID?

- Do we really need to send the state ID?

- Fixed width messages?

- A possible message from state to central system:

$$25 \quad AAAAPPPPPBBBBBWWWSSOOOOOAAA$$

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Outline

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (basic idea)

### RLE Algorithm

The idea is counting the times that each character appears consecutively. For example, for a string:

$$S = aaaabbbbbbbbaaaaabbbbbbccccccbb$$

its compressed representation will be:

$$\tilde{S} = a4b8a5b6c5b2$$

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (algorithm v1.0)

```
function char * compress(const char *input) begin
    char *str ← input;
    char *output ← new char;
    int length ← 0;
    while *str ≠ 0 do
        char x ← *str;
        push-back(output,x);
        int k ← 1;
        while x = *(++str) do k++;
        push-back(output,to-alpha(k));
        length ← length + k + 1;
    end
    return strlen(output) < length ? output : input;
end
```

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (inconvenients)

- What about decompression?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (inconvenients)

- What about decompression? Very simple, duh!

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (inconvenients)

- What about decompression? Very simple, duh!
- What happens if we got a compressed (and ambiguos) string like following?

$$\tilde{S} = a315b3$$

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (inconvenients)

- What about decompression? Very simple, duh!
- What happens if we got a compressed (and ambiguos) string like following?

$$\tilde{S} = a315b3$$

- Maybe, original string was like:

$$S = aaa11111bbbb$$

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

**RLE Algorithm**
RLE (improved)
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (inconvenients)

- What about decompression? Very simple, duh!
- What happens if we got a compressed (and ambiguos) string like following?

$$\tilde{S} = a315b3$$

- Maybe, original string was like:

$$S = aaa11111bbbb$$

- Or maybe was:

$$S = aaa \cdots aaabbbb$$

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (inconvenients)

- What about decompression? Very simple, duh!
- What happens if we got a compressed (and ambiguos) string like following?

$$\tilde{S} = a315b3$$

- Maybe, original string was like:

$$S = aaa11111bbbb$$

- Or maybe was:

$$S = aaa \cdots aaabbbb$$

- Is this algorithm effective with XML or JSON?

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Outline

Ulises Tirado Zatarain        Data Compression & IoT

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- Maybe, we can use a separator/delimiter character?

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- Maybe, we can use a separator/delimiter character?
  - What character can we use to «,»,«|»,...?

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- Maybe, we can use a separator/delimiter character?
  - What character can we use to «,»,«|»,...?
  - What if this character is in the original message?

$$S = ||||,,, \Rightarrow \tilde{S} = |4,,3 \text{ or } \tilde{S} = |4|,3$$

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- Maybe, we can use a separator/delimiter character?
  - What character can we use to «,»,«|»,...?
  - What if this character is in the original message?

  $$S = ||||,,, \Rightarrow \tilde{S} = |4,,3 \text{ or } \tilde{S} = |4|,3$$

- What if we only have two kinds of characters?

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- Maybe, we can use a separator/delimiter character?
  - What character can we use to «,»,«|»,...?
  - What if this character is in the original message?

$$S = ||||,,, \Rightarrow \tilde{S} = |4,,3 \text{ or } \tilde{S} = |4|,3$$

- What if we only have two kinds of characters?
  - Thinking in binary :)

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- Maybe, we can use a separator/delimiter character?
  - What character can we use to «,»,«|»,...?
  - What if this character is in the original message?
  
  $$S = |||||,,, \Rightarrow \tilde{S} = |4,,3 \text{ or } \tilde{S} = |4|,3$$

- What if we only have two kinds of characters?
  - Thinking in binary :) The Beattles - Let it «bit»! XD

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- Maybe, we can use a separator/delimiter character?
  - What character can we use to «,»,«|»,...?
  - What if this character is in the original message?

$$S = |||||,,, \Rightarrow \tilde{S} = |4,,3 \text{ or } \tilde{S} = |4|,3$$

- What if we only have two kinds of characters?
  - Thinking in binary :) The Beattles - Let it «bit»! XD

$$S = "e >< v" = \quad 00100101 \quad 00111110 \quad 00111100 \quad 01110110$$
$$\tilde{S} = 0100010100010010010101010101110001101100101001$$

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- Maybe, we can use a separator/delimiter character?
  - What character can we use to «,»,«|»,...?
  - What if this character is in the original message?

$$S = \||||,,, \Rightarrow \tilde{S} = |4,,3 \text{ or } \tilde{S} = |4|,3$$

- What if we only have two kinds of characters?
  - Thinking in binary :) The Beatles - Let it «bit»! XD

$$S = "e > < v" = \quad 00100101 \quad 00111110 \quad 00111100 \quad 01110110$$
$$\tilde{S} = 0100010100010010010101010101110001101100101 0001$$

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- Maybe, we can use a separator/delimiter character?
  - What character can we use to «,»,«|»,...?
  - What if this character is in the original message?

$$S = |||||,, \Rightarrow \tilde{S} = |4,,3 \text{ or } \tilde{S} = |4|,3$$

- What if we only have two kinds of characters?
  - Thinking in binary :) The Beattles - Let it «bit»! XD

$$S = "e >< v" = \begin{array}{cccc} 00100101 & 00111110 & 00111100 & 01110110 \end{array}$$
$$\tilde{S} = 010001010001001001010101011110001101100101001$$

  **FAIL!!**

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- Maybe, we can use a separator/delimiter character?
  - What character can we use to «,»,«|»,...?
  - What if this character is in the original message?

$$S = ||||,,, \Rightarrow \tilde{S} = |4,,3 \text{ or } \tilde{S} = |4|,3$$

- What if we only have two kinds of characters?
  - Thinking in binary :) The Beattles - Let it «bit»! XD

$$S = "e >< v" = \quad 00100101 \quad 00111110 \quad 00111100 \quad 01110110$$
$$\tilde{S} = 0100010100010010010101010111000110110010100001$$

  **FAIL!!**

- What if we limit the repetitions? (i.e. MAX 9 repetitions)

$$S = aaaaaaaaaaaabbb$$
$$\tilde{S} = a9a3b$$

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- What if the number of repetitions always are represented by only one byte?

  - What is the maximum repetition for a single character in this representation?
  - How can detect corrupted data from compressed message?

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- What if the number of repetitions always are represented by only one byte?

$$\tilde{S} = a!b\#$$

This compressed string represents 33 a's followed by 35 b's. Why?

  - What is the maximum repetition for a single character in this representation?
  - How can detect corrupted data from compressed message?

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- What if the number of repetitions always are represented by only one byte?

$$\tilde{S} = a!b\#$$

This compressed string represents 33 $a$'s followed by 35 $b$'s. Why?

- The ASCII code of ! is 33 and the code of # is 35. Then, in this representation the 2k-th characters are basic information and $(2k+1)$-th characters are number of repetitions, where $k = 0, 1, 2, \ldots, \frac{n}{2}$.
- What is the maximum repetition for a single character in this representation?
- How can detect corrupted data from compressed message?

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- What if the number of repetitions always are represented by only one byte?

$$\tilde{S} = a!b\#$$

  This compressed string represents 33 $a$'s followed by 35 $b$'s. Why?

  - The ASCII code of ! is 33 and the code of # is 35. Then, in this representation the 2k-th characters are basic information and $(2k+1)$-th characters are number of repetitions, where $k = 0, 1, 2, \ldots, \frac{n}{2}$.
  - What is the maximum repetition for a single character in this representation?
  - How can detect corrupted data from compressed message?

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

- What if the number of repetitions always are represented by only one byte?

$$\tilde{S} = a!b\#$$

This compressed string represents 33 $a$'s followed by 35 $b$'s. Why?

  - The ASCII code of ! is 33 and the code of # is 35. Then, in this representation the 2k-th characters are basic information and $(2k+1)$-th characters are number of repetitions, where $k = 0, 1, 2, \ldots, \frac{n}{2}$.
  - What is the maximum repetition for a single character in this representation?
  - How can detect corrupted data from compressed message?

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
**RLE (improved)**
RLE (decompression)
RLE (PRO)

# Run-Length Encoding (improved v2.0)

```
function char * compress(const char *input)begin
    const char MAX ←~ 0;
    char *str ← input;
    char *output ← new char;
    int length ← 0;
    while *str ≠ 0 do
        char x ← *str;
        push-back(output,x);
        char k ← 1;
        while x = *(++str) && k < MAX do k++;
        push-back(output,k);
        length ← length+k+1;
    end
    return strlen(output) < length ? output : input;
end
```

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
**RLE (decompression)**
RLE (PRO)

# Outline

Introduction
RLE Algorithm
**Basic algorithms (Loseless)** RLE (improved)
Advanced algorithms (Loseless) **RLE (decompression)**
Advanced algorithms (Lossy) RLE (PRO)

# Run-Length Encoding (decompression)

```
function char * decompress(const char *input) begin
    char *c ← input;
    char *output ← new char;
    while *c ≠ 0 do
        char n ← *(c + 1);
        for i = 1 to n do  push-back(output,*c) ;
        c ← c + 2;
    end
    return  output;
end
```

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
**RLE (PRO)**

# Outline

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
**RLE (PRO)**

# Run-Length Encoding (improved v3.0)

Maybe we can do better:

- We can analize input message to see what unique characters are contained in the mesage.

- Then, we can construct a reduced alphabet. For example, $\Sigma = \{A, B, O, P, S, W\}$.

- In this case only need 3 bits to represent characters.

- We can use 5 bits to store repetitions (so, we will have MAX 31 repetitions by character).

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
**RLE (PRO)**

# Run-Length Encoding (improved v3.0)

Maybe we can do better:

- We can analize input message to see what unique characters are contained in the mesage.

- Then, we can construct a reduced alphabet. For example, $\Sigma = \{A, B, O, P, S, W\}$.

- In this case only need 3 bits to represent characters.

- We can use 5 bits to store repetitions (so, we will have MAX 31 repetitions by character).

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
**RLE (PRO)**

# Run-Length Encoding (improved v3.0)

Maybe we can do better:

- We can analize input message to see what unique characters are contained in the mesage.

- Then, we can construct a reduced alphabet. For example, $\Sigma = \{A, B, O, P, S, W\}$.

- In this case only need 3 bits to represent characters.

- We can use 5 bits to store repetitions (so, we will have MAX 31 repetitions by character).

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
**RLE (PRO)**

# Run-Length Encoding (improved v3.0)

Maybe we can do better:

- We can analize input message to see what unique characters are contained in the mesage.
- Then, we can construct a reduced alphabet. For example, $\Sigma = \{A, B, O, P, S, W\}$.
- In this case only need 3 bits to represent characters.
- We can use 5 bits to store repetitions (so, we will have MAX 31 repetitions by character).

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
**RLE (PRO)**

# Run-Length Encoding (improved v3.0)

Maybe we can do better:

- We can analize input message to see what unique characters are contained in the mesage.
- Then, we can construct a reduced alphabet. For example, $\Sigma = \{A, B, O, P, S, W\}$.
- In this case only need 3 bits to represent characters.
- We can use 5 bits to store repetitions (so, we will have MAX 31 repetitions by character).

Introduction
**Basic algorithms (Loseless)**
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

RLE Algorithm
RLE (improved)
RLE (decompression)
**RLE (PRO)**

# Run-Length Encoding (improved v4.0 PRO)

Moreover, what happens if we can detect whole words?

$$S = abcdabcdabcdbdbdbdbd$$

$$\tilde{S} = [abcd]\,3\,[bd]\,4$$

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

# Outline

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

# Probability basics

## Experiments and Events

An event is a set of posible results in an experiment execution. For example, taking a card from a deck or rolling a dice. We can denote all posible results with $\Omega$ and an event with uppercase letter such that $A \subseteq \Omega$ or $A \in 2^{\Omega}$, then $A$ is a subset of $\Omega$.

## Probability

Is an indicator that describes the frecuency of an event in one universal set of posibilities. Daily, we express that indicator as a percentage value or value between $0$ and $1$. Then we can define the probability as:

$$p : 2^{\Omega} \mapsto [0, 1]$$

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

## Probability basics

- So, when $\Omega$ is a discrete and finite set, then:

$$p(A) = \frac{\#A}{\#\Omega}$$

- We name this as uniform distribution or counting distribution.
- However, counting elements in $A$ and $\Omega$ isn't always trivial. Maybe we need to use operations like factorial, combinations, permutations, etcetera.

Ulises Tirado Zatarain        Data Compression & IoT

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

## Probability basics

- Key pressing random letter of the english keyboard such that the letter be a vowel.
  - Let $A = \{a, e, i, o, u\}$ and $\Omega = \{a, ..., z\}$, then $p(A) = \frac{5}{26}$.
- Rolling a dice such that the result be greater than 2.
  - Let $A = \{3, 4, 5, 6\}$ and $\Omega = \{1, 2, 3, 4, 5, 6\}$, then $p(A) = \frac{2}{3}$.,
- Taking a card from a deck such that getting a red card.
  - Let
    $A = \{A\heartsuit, 2\heartsuit, ..., 10\heartsuit, J\heartsuit, Q\heartsuit, K\heartsuit, A\diamondsuit, 2\diamondsuit, ..., 10\diamondsuit, J\diamondsuit, Q\diamondsuit, K\diamondsuit\}$
    and $\Omega =$
    $$\left\{ \begin{array}{ll} A\heartsuit, 2\heartsuit, ..., 10\heartsuit, J\heartsuit, Q\heartsuit, K\heartsuit, & A\diamondsuit, 2\diamondsuit, ..., 10\diamondsuit, J\diamondsuit, Q\diamondsuit, K\diamondsuit, \\ A\spadesuit, 2\spadesuit, ..., 10\spadesuit, J\spadesuit, Q\spadesuit, K\spadesuit, & A\clubsuit, 2\clubsuit, ..., 10\clubsuit, J\clubsuit, Q\clubsuit, K\clubsuit \end{array} \right\},$$
    then $p(A) = \frac{1}{2}$.

Introduction
Basic algorithms (Loseless)
**Advanced algorithms (Loseless)**
Advanced algorithms (Lossy)

Probability Review
**Data Structures Review**
Huffman Code

# Outline

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

# Priority Queue

## Wikipedia

- In computer science, a priority queue is an abstract data type which is like a regular queue or stack data structure, but where additionally each element has a "priority" associated with it. In a priority queue, an element with high priority is served before an element with low priority.

- While priority queues are often implemented with heaps, they are conceptually distinct from heaps. A priority queue is an abstract concept like "a list" or "a map";

Ulises Tirado Zatarain          Data Compression & IoT

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

# C++ STL Priority Queue

## std::priority_queue

template $<$class T, class Container $=$ vector$<$T$>$,
class Compare $=$ less$<$typename Container::value_type$>>$ class
priority_queue;

**Priority queue**
Priority queues are a type of container adaptors, specifically designed such that its first element is always the greatest of the elements it contains.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

The container shall be accessible through random access iterators and support the following operations:

- empty
- size
- front
- push_back
- pop_front

The standard container classes vector and deque fulfill these requirements. By default, if no container class is specified for a particular priority_queue class instantiation, the standard container vector is used.

Introduction
Basic algorithms (Loseless)
**Advanced algorithms (Loseless)**
Advanced algorithms (Lossy)

Probability Review
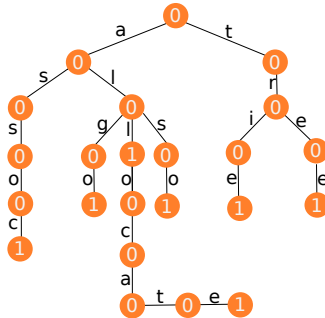**Data Structures Review**
Huffman Code

# Tries

- Let a **word** be a single string and let **dictionary** be a large set of words.
- The **set<string>** and the **hash tables** can only find in a dictionary words that match exactly with the single word that we are finding.
- **Trie** is a tree type data structure that allows to represent a dictionary.
  - We can insert and find strings in $\mathcal{O}(L)$.
  - We can perform incremental search.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

# Tries

- The word **trie** is an infix of the word "re**trie**val" because the trie can find a single word in a dictionary with only a prefix of the word.
- The trie is a tree where each vertex represents a single word or a prefix.
  - The root represents an empty string $\varepsilon$.
  - A vertex that are $k$ edges of distance of the root have an associated prefix of length $k$.
  - Let $v$ and $w$ be two vertexes of the trie, and assume that $v$ is a direct father of $w$, then $v$ must have an associated prefix of $w$.
  - Deterministic acyclic finite state automaton.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

## Tries: Example

- The following trie stores the words: "algo", "assoc", "all", "allocate", "also", "tree" and "trie".



- Note that every vertex of the tree does not store entire prefixes or entire words.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

## How to represent tries?

- The most simple way to represent a trie is with an struct like following:

  **struct** trie $<$typename Type $=$ int$>$ **begin**
  > Type data;
  > **struct** trie *edge[ALPHABET_SIZE];

  **end**

- For the english alphabet, we can store the $'a'$-edge in $trie::edge[0]$, $'b'$-edge in $trie::edge[1]$, $'c'$-edge in $trie::edge[2]$ and so on until $'z'$-edge in $trie::edge[25]$.

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

# How to add a word to dictionary?

- We can add a word $w$ as following:

```
function add-word(struct trie *t , char *w) begin
    if is-empty(w) then:
        │  t → data ⟵ t → data + 1;
    else:
        │  if is-null(t → edge [*w]) then:
        │  │  t → edge [*w] ⟵ new struct trie;
        │  end
        │  add-word(t → edge [*w], w + 1);
    end
end
```

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

# How to find a word in dictionary?

- To find a word $w$, we can perform following algorithm:

  **function** find-word(**struct** trie $*$t , **char** $*w$)**begin**

      **if** is-null(t) **then:**

          | **return** $0$;

      **end**

      **if** is-empty($w$) **then:**

          | **return** t $\rightarrow data$;

      **end**

      **return** find-word(t $\rightarrow edge[*w], w+1$);

  **end**

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

# Outline

Introduction
Basic algorithms (Loseless)
Advanced algorithms (Loseless)
Advanced algorithms (Lossy)

Probability Review
Data Structures Review
Huffman Code

# Huffman Encoding

- The idea is...

## Statistical

- Statistical Pattern Recognition
  - Principal Component Analysis
- Digital Signal Processing
  - Filtering
- Image compression
  - Grayscale images
  - Color images

# References I

- Standford University
- HackerRank
- Code Forces
- Code Chef
- Wikipedia