# Lecture 7

**Asymmetric Numeral Systems**

# Announcments

→ HW1 due TODAY , 11.59 pm !

→ HW2 released by tonight / tomorrow.
  Due 11/1 , Wed, 11.59 pm

→ Project details released this week.
  Proposal due 11/1 , Wed, 11.59 pm

→ Coding : use coding gradescope link

# Quiz Q1

Consider a Bernoulli random variable (we have been seeing them a lot in the quizzes so far :)) - $Ber(p)$ where $P(A) = p$ and $P(B) = (1 - p)$. Consider sequence of symbols AB and BA , to be encoded using Arithmetic Coding (assume idealized version where the codelength is equal to $\log_2\left(\frac{1}{P(X^n)}\right)$).
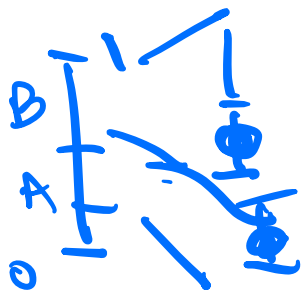
## Q1.1

AB and BA have the same code-length.

(✓) True

( ) False

## Q1.2

AB and BA have the same output codeword.

( ) True

(✓) False

$$P(AB) = P(1-P)$$

$$P(BA) = (1-P)P$$

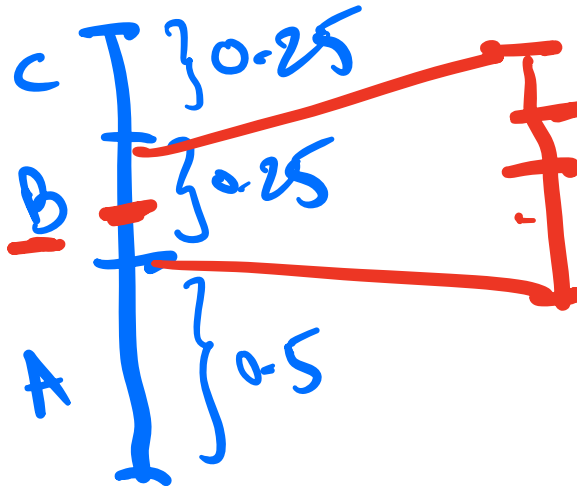$$L(x^n) \lessgtr \lceil \log(1/P(x^n)) \rceil + 1$$

# Q1.2

Assume a probability distribution over symbols $\mathcal{X} = \{A, B, C\}$ with respective probabilities $p(A) = 0.5, p(B) = 0.25, p(C) = 0.25$. An arithmetic decoder receives as input bitstream `100111` for an input of length 3. What is the decoded sequence?

## Ans

b `100111` = 0.609375

Decoded codeword = `BAC`

# Slides credit - Kedar Tatwawadi

# RECAP -> Arithmetic coding

1. Given *any* distribution $P$, achieves *optimal* compression. Thus, Arithmetic coding allows for `model` and `entropy coding` separation.

2. Can work very well with changing distribution $P$.
   i.e. adaptive algorithms work well with Arithmetic coding

# RECAP -> Arithmetic coding in practice

Lots of Variants of Arithmetic coding; mainly come from how they implement the rescaling.

1. **Arithmetic coding:** Bit-based rescaling -> keeping a count of the mid-ranges etc.
   SCL Arithmetic coding

2. **Range Coding** Byte (8-bit based rescaling), word-based rescaling ->
   SCL range coding

# RECAP -> Arithmetic/Range coders in practice

Used almost everywhere! (either as Range coder or Arithmetic coding)

1. JPEG2000, BPG, H265, H266, VP8

2. CMIX, `tensorflow-compress`, NNCP

# RECAP -> Arithmetic/Range coders in practice

Although Arithmetic coding algorithms are quite efficient, they are not fast enough!
(especially when compared with Huffman coding)

| Codec | Encode speed | Decode speed | compression |
|---|---|---|---|
| Huffman coding | 252 Mb/s | 300 Mb/s | 1.66 |
| Arithmetic coding | 120 Mb/s | 69 Mb/s | 1.24 |

NOTE -> Speed numbers from: Charles Bloom's blog

# RECAP -> Arithmetic/Range coders in practice

| Codec | Encode speed | Decode speed | compression |
|---|---|---|---|
| Huffman coding | 252 Mb/s | 300 Mb/s | 1.66 |
| Arithmetic coding | 120 Mb/s | 69 Mb/s | 1.24 |
| rANS | 76 Mb/s | 140 Mb/s | 1.24 |
| tANS | 163 Mb/s | 284 Mb/s | 1.25 |

NOTE -> Speed numbers from: Charles Bloom's blog

# ANS: Asymmetric Numeral System

- `rANS` : *range-ANS*: drop-in replacement for Arithmetic coding, but faster
- `tANS` : *table-ANS*: drop-in replacement for Huffman coding, but better in terms of compression (and similar speed)

# Why *"Asymmetric Numeral System"*?

# Why *"Asymmetric Numeral System"*?

Lets assume inputs are digits `[0,9]`

```
data_input = [3,2,4,1,5]
```

Quiz-1: Form a single number `x` (*state*) representing `data_input` ?

# Why *"Asymmetric Numeral System"*?

Lets assume inputs are digits `[0,9]`

```
data_input = [3,2,4,1,5]
```

Quiz-1: Form a single number `x` (*state*) representing `data_input` ?

*Ans:* `[3,2,4,1,5] -> 32415`

$$3$$

$$(3,2) = 3 \times 10 + 2$$
$$= 32$$

$$(3,2,4) = 3 \times 100$$
$$+ 2 \times 10$$
$$+ 4$$

$$= (32) \times 10$$
$$+ 4$$

$$\vdots$$

# Symmetric Numeral System: *Encoding*

```python
# given
data_input = [3,2,4,1,5]

## "encoding" process
x = 0   # <-- initial state
x = x*10 + 3 #x = 3
x = x*10 + 2 #x = 32
x = x*10 + 4 #x = 324
x = x*10 + 1 #x = 3241
x = x*10 + 5 # x = 32415 <- final state
```

$$s = 32145$$

$$(s//10, s\%10)$$

$$= (3214, 5)$$

Sprev    symbol

# Symmetric Numeral System: *Decoding*

Quiz-2: Given the *state* `x`, how can you retrieve the `data_input`?

```
symbols = []
x = 32145 # <- final state
n = 5
```

# Symmetric Numeral System: *Decoding*

Quiz-2: Given the *state* `x` , and number of elements `n=5` how can you retrieve the `data_input` ?

```
symbols = []
x = 32145 # <- final state
n = 5

# repeat n=5 times
s, x = x%10, x//10   # s,x = 5, 3214
s, x = x%10, x//10   # s,x = 4, 321
s, x = x%10, x//10   # s,x = 1, 32
s, x = x%10, x//10   # s,x = 2, 3
s, x = x%10, x//10   # s,x = 3, 0
```

# Symmetric Numeral System: *Encoding, Decoding*

`encode_step` , `decode_step` are exact inverses of each other

```python
def encode_step(x_prev,s):
    x = x_prev*10 + s
    return x

def decode_step(x)
    s = x%10 # decode symbol
    x_prev = x//10 # retrieve the previous state
    return (s,x_prev)
```

# Symmetric Numeral System

```
symbols = []
x = 32145 # <- final state
n = 5
```

```
# Encoding
state x: 0 -> 3 -> 32 -> 321 -> 3214 -> 32145

# decoding
state x: 0 <- 3 <- 32 <- 321 -> 3214 <- 32145
```

# Symmetric Numeral System: Full *Encoding*

Transmit the final state `x` in binary using `ceil(log_2[x+1]))` .
(eg. `32145 = b111110110010001` ).

```
## Encoding
def encode_step(x_prev,s):
    x = x_prev*10 + s
    return x

def encode(data_input):
    x = 0 # initial state
    for s in data_input:
        x = encode_step(x,s)

    return to_binary(x) # takes ~log_2(x) bits
```

$$s \in \{0, 1, \ldots, 9\}$$

$$\underline{x} \longmapsto \underline{10 * x}$$

$$\sim \log_2(10)$$

*per symbol*

# Symmetric Numeral System: Full *Decoding*

```python
## Decoding
def decode_step(x):
    s = x%10 # <- decode symbol
    x_prev = x//10 # <- retrieve the previous state
    return (s,x_prev)

def decode(bits, num_symbols):
    x = to_uint(bits) # convert bits to final state

    # main decoding loop
    symbols = []
    for _ in range(num_symbols):
        s, x = decode_step(x)
        symbols.append(s)
    return reverse(symbols) # need to reverse to get original sequence
```

# Symmetric Numeral System: Compression performance

- `x <- 10*x_prev + s` , i.e `x ~ x_prev*10`

- Per symbol we are using approximately:

$$\log_2(10) = \log_2 \frac{1}{0.1} \quad = \log_2\left(1/p\right)$$

- Our compressor is "optimal", only if all our symbol `{0, ..., 9}` have equal probability of `0.1` .

$$0 : p_0$$
$$1 \rightarrow p_1$$

# Asymmetric Numeral Systems:

Quiz 2: If that digits `[0,9]` have different probabilities, how can we modify our algorithm to achieve better compression?
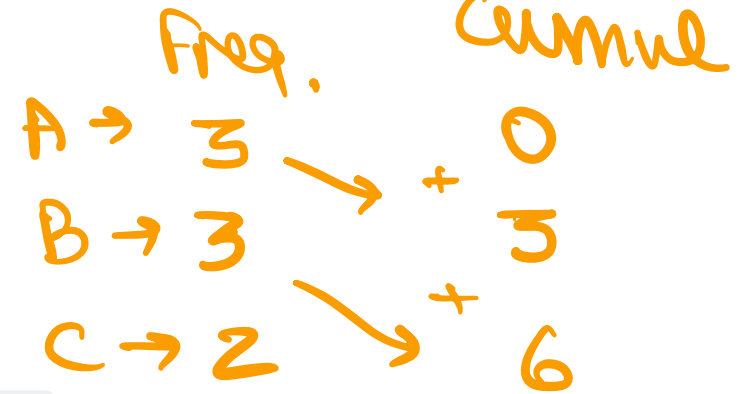
$$x \approx x_{prev} * 1/p_s$$

$$\approx \log 1/p_s$$

# Asymmetric Numeral Systems:

If that digits `[0,9]` have different probabilities, how can we modify our algorithm to achieve better compression?

*Ans:* Instead of scaling `x ~ 10*x_prev`, we want to scale it as `x ~ (1/prob[s])*x_prev`

# rANS: Notation



- represent probabilities in terms of integer `frequencies`

```
prob[s] = freq[s]/M
```

- For example:

```
# prob in terms of integers
alphabet  -> {0,1,2}, prob -> [3/8, 3/8, 2/8]
freq -> [3,3,2], M = 8
cumul -> [0,3,6]
```

# rANS: Encoding

```python
def rans_base_encode_step(x_prev,s):
    # TODO: add details here
    return x


def encode(data_input):
    x = 0 # initial state
    for s in data_input:
        x = rans_base_encode_step(x,s)

    return to_binary(x)
```

$$x = 10 * x + s$$
return x

Quiz 4: What is the `rans_base_encode_step` ?
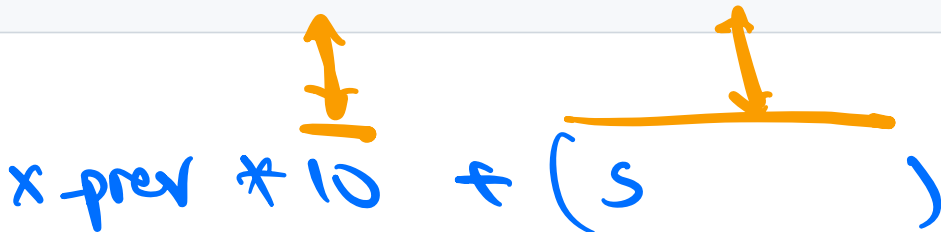
(*HINT*: `x ~ (1/prob[s])*x_prev` )

# rANS: Encoding

```python
def rans_base_encode_step(x_prev,s):
    x = (x_prev//freq[s])*M + cumul[s] + x_prev%freq[s]
    return x

def encode(data_input):
    x = 0 # initial state
    for s in data_input:
        x = rans_base_encode_step(x,s)

    return to_binary(x)
```

(*HINT*: `x ~ (1/prob[s])*x_prev` )

# rANS: Encoding

```python
def rans_base_encode_step(x_prev,s):
    x = (x_prev//freq[s])*M + cumul[s] + x_prev%freq[s]
    return x
```

$$x\_prev * 10 \quad + \quad ( s \qquad )$$

# rANS: Encoding

$x = 0$

$s = 1, \quad x = (0 // 3) * 8 + 3 +$
$+ (0 \% 3)$
$= 3$

```python
def rans_base_encode_step(x_prev,s):
    x = (x_prev//freq[s])*M + cumul[s] + x_prev%freq[s]
    return x
```

*Example-1*

$s = 0, \quad x = (3 // 3) * 8 + 0$
$+ (3 \% 3)$
$= 8$

```
alphabet  -> {0,1,2}
freq -> [3,3,2], M = 8, cumul -> [0,3,6]
```

Quiz 4: What is the final encoding for `data_input = [1,0,2,1]` ?

# rANS: Encoding

```
def rans_base_encode_step(x_prev,s):
    x = (x_prev//freq[s])*M + cumul[s] + x_prev%freq[s]
    return x
```

*Example-1:*

```
freq -> [3,3,2], M = 8, cumul -> [0,3,6]
```

```
step 0: x = 0
step 1: s = 1 -> x = (0//3)*8 + 3 + (0%3) = 3
step 2: s = 0 -> x = (3//3)*8 + 0 + (3%3) = 8
step 3: s = 2 -> x = (8//2)*8 + 6 + (8%2) = 38
step 4: s = 1 -> x = (38//3)*8 + 3 + (38%3) = 101
```

# rANS: Encoding

```python
def rans_base_encode_step(x_prev,s):
    x = (x_prev//freq[s])*M + cumul[s] + x_prev%freq[s]
    return x
```

*Example-2*

$$x = x_{prev} * 10 + s$$

```
alphabet = {0,1,2,3,4,5,6,7,8,9}
freq[s] = 1, M = 10
prob[s] = 1/10
```

Quiz 5: How does the `rans_base_encode_step` function look for the input above?

# rANS: Encoding

$$slot = cumul(s) + x\_prev \% freq(s)$$

$$< freq(s)$$

$$\geq 0$$

```
x = (x_prev//freq[s])*M + cumul[s] + x_prev%freq[s]
```

```
def rans_base_encode_step(x_prev,s):
    # find block_id, slot
    block_id = x_prev//freq[s]
    slot = cumul[s] + (x_prev % freq[s])

    # compute next state x
    x = block_id*M + slot
    return x
```

① ② 

$$0 \leq cumul(s) \leq slot < cumul(s+1)$$
$$< M$$

Quiz 6: What can you say about the `slot` ?

$$x = block\_id * M + \underline{slot}$$
$$\epsilon [0, M)$$

# rANS: Encoding

```
x = (x_prev//freq[s])*M + cumul[s] + x_prev%freq[s]
```

```python
def rans_base_encode_step(x_prev,s):
    # find block_id, slot
    block_id = x_prev//freq[s]
    slot = cumul[s] + (x_prev % freq[s])

    # compute next state x
    x = block_id*M + slot
    return x
```

Quiz 6: What can you say about the `slot` ?

*Ans*: `0 <= slot < M`

# rANS: Decoding

```python
def rans_base_encode_step(x_prev,s):
    block_id = x_prev//freq[s]
    slot = cumul[s] + (x_prev % freq[s])
    x = block_id*M + slot
    return x

def rans_base_decode_step(x):
    # TODO -> fill in the decoder here
    return (s,x_prev)
```

# rANS: Decoding

```python
def rans_base_encode_step(x_prev,s):
    block_id = x_prev//freq[s]
    slot = cumul[s] + (x_prev % freq[s])
    x = block_id*M + slot
    return x

def rans_base_decode_step(x):
    # Step I: find block_id, slot
    block_id = ?
    slot = ?

    ...
    return (s,x_prev)
```

STEP-I: Decode `block_id, slot`

# rANS: Decoding

$$slot = cumul[s] + (x\_prev \% freq[s])$$

```
def rans_base_encode_step(x_prev,s):
    block_id = x_prev//freq[s]
    slot = cumul[s] + (x_prev % freq[s])
    x = block_id*M + slot
    return x

def rans_base_decode_step(x):
    block_id, slot = x//M, x%M # Step I

    s = ? # Step II: decode symbol s
    ...
    return (s,x_prev)
```

$$0 \le \quad slot \quad < M$$

$$cumul[s] \le \quad < cumul[s+1]$$

$cumul[s]$

STEP-II: Decode symbol s

# rANS: Decoding

```python
def rans_base_encode_step(x_prev,s):
    block_id = x_prev//freq[s]
    slot = cumul[s] + (x_prev % freq[s])
    x = block_id*M + slot
    return x

def rans_base_decode_step(x):
    block_id, slot = x//M, x%M # Step I

    s = ? # Step II: decode symbol s
    ...
    return (s,x_prev)
```

*HINT*: `cumul[s] <= slot < cumul[s+1]`

# rANS: Decoding

```python
def rans_base_encode_step(x_prev,s):
    block_id = x_prev//freq[s]
    slot = cumul[s] + (x_prev % freq[s])
    x = block_id*M + slot
    return x

def rans_base_decode_step(x):
    block_id, slot = x//M, x%M # Step I
    s = find_bin(cumul, slot) # Step II: decode symbol s
    ...
    return (s,x_prev)
```

*HINT*:  `cumul[s] <= slot < cumul[s+1]`

s ✓

block_id
$= x\_prev // freq[s]$

slot
$= cumul[s]$
$+ (. \sim \% ..)$

# rANS: Decoding

```python
def rans_base_encode_step(x_prev,s):
    block_id = x_prev//freq[s]
    slot = cumul[s] + (x_prev % freq[s])
    x = block_id*M + slot
    return x

def rans_base_decode_step(x):
    block_id, slot = x//M, x%M # Step I
    s = find_bin(cumul, slot) # Step II: decode symbol s

    x_prev = ? # Step-III
    return (s,x_prev)
```

STEP-III: Decode prev state `x_prev`

$$x_{\cdot} = \left( x_{prev} \mathbin{//} freq[s] \right) + cumul[s]$$
$$+ \ x_{prev} \, \% \, freq[s] \ \checkmark$$

# rANS: Decoding

```python
def rans_base_encode_step(x_prev,s):
    block_id = x_prev//freq[s]
    slot = cumul[s] + (x_prev % freq[s])
    x = block_id*M + slot
    return x

def rans_base_decode_step(x):
    block_id, slot = x//M, x%M # Step I
    s = find_bin(cumul, slot) # Step II: decode symbol s
    x_prev = block_id*freq[s] + slot - cumul[s] # Step-III
    return (s,x_prev)
```

# rANS: Decoding example

x-final = 101

block_id, slot = 12, 5

s = 1 (cumul[s] = 3)

```python
def rans_base_decode_step(x):
    block_id, slot = x//M, x%M # Step I
    s = find_bin(cumul, slot) # Step II: decode symbol s
    x_prev = block_id*freq[s] + slot - cumul[s] # Step-III
    return (s,x_prev)
```

x-prev = 12 * 3 + 5 - 3
= 38

*Example-1*:

5

```
freq -> [3,3,2], M = 8, cumul -> [0,3,6]
```

(1, 38)

```
# Decode
bitarray = b1100101 -> x_final = 101
n = 4
```

# rANS: Encoder, Decoder

Interactive snippet: https://kedartatwawadi.github.io/post--ANS/

# rANS: Full Encoder

```python
def rans_base_encode_step(x_prev,s):
    block_id = x_prev//freq[s]
    slot = cumul[s] + (x_prev % freq[s])
    x = block_id*M + slot
    return x


def encode(data_input):
    x = 0 # initial state
    for s in data_input:
        x = rans_base_encode_step(x,s)
    return to_binary(x)
```

# rANS: Full Decoder

```python
def rans_base_decode_step(x):
    block_id, slot = x//M, x%M # Step I
    s = find_bin(cumul, slot) # Step II: decode symbol s
    x_prev = block_id*freq[s] + slot - cumul[s] # Step-III
    return (s,x_prev)

def decode(bits, num_symbols):
    x = to_uint(bits) # convert bits to final state

    # main decoding loop
    decoded_symbols = []
    for _ in range(num_symbols):
        s, x = rans_base_decode_step(x)
        decoded_symbols.append(s)
    return reverse(decoded_symbols) # need to reverse to get original sequence
```

$$x \approx \left( \frac{1}{p(s_1)} \times \frac{1}{p(s_2)|s_1} \cdots \right)$$

# rANS Compression performance:

- state increases as: $x \approx x_{prev} \cdot \frac{M}{freq[s]} = x_{prev} \cdot \frac{1}{P(s)}$

- Final state $x_{final}$ represented using $\approx \log_2(x_{final})$ bits.

Quiz 8: What is the approximate encode length for input $s^n$?

# rANS Compression performance:

- state increases as: $x \approx x_{prev} \cdot \frac{M}{freq[s]} = x_{prev} \cdot \frac{1}{P(s)}$

- Final state $x_{final}$ represented using $\approx \log_2(x_{final})$ bits.

Quiz 8: What is the approximate encode length for input $s^n$?

$$L(s^n) \approx \log_2 \frac{1}{P(s^n)}$$

i.e. `rANS` is *optimal*!

# Connection with Bits-Back Coding

- Another way to think about rANS and why it achieves optimal compression.

- Connection with $H(X) = H(X, Z) - H(Z|X)$

- More on this soon:

  - See HW2, Q4 on bits-back coding

  - More on conditional entropy in next few lectures

# rANS vs Arithmetic coding

- **Optimal Compression**: Compression performance *optimal* and similar to Arithmetic coding in practice

- **Reverse decoding**: Decoding is in the *reverse* order unlike Arithmetic coding, which can be a bit of a problem for Adaptive schemes

- **Simpler encoding/decoding**: Encoding, Decoding operations are a bit simpler, making it easier to modify and optimize for speed

# rANS in practice

Quiz 9: What is the practical problem with our rANS Encoding/Decoding?

**RANS ENCODING EXAMPLE**

Symbol Counts, $\mathcal{F}$ [3,3,2]
Input Symbol String: [1,0,2,1,0,2,2,1,0,1, 2,2,2,2]
[Try it]

| Input | State |
|-------|----------|
| 1 | 3 |
| 0 | 8 |
| 2 | 38 |
| 1 | 101 |
| 0 | 266 |
| 2 | 1070 |
| 2 | 4286 |
| 1 | 11429 |
| 0 | 30474 |
| 1 | 81267 |
| 2 | 325071 |
| 2 | 1300287 |
| 2 | 5201151 |
| 2 | 20804607 |

# rANS in practice

Quiz 9: What is the practical problem with our rANS Encoding/Decoding?

- The state increases as: $x \approx x_{prev} \cdot \frac{M}{freq[s]} = x_{prev} \cdot \frac{1}{P(s)}$

- After `~30-40` symbols we are going to need more than `64` bits

**Solution**: Restrict the state `x` to lie in an interval `[L,H]`

# Streaming rANS: Encoding

```python
def rans_base_encode_step(x_shrunk, s):
    ...

def shrink_state(x_prev,s):
    ...
    return x_shrunk, out_bits

def encode_step(x_prev,s):
    # shrink state x before calling base encode
    x_shrunk, out_bits = shrink_state(x, s)

    # perform the base encoding step
    x = rans_base_encode_step(x_shrunk,s)
    assert x in [L,H]
    return x, out_bits
```

# Streaming rANS: shrink state

```python
def encode_step(x_prev,s):
    # shrink state x before calling base encode
    x_shrunk, out_bits = shrink_state(x, s)

    # perform the base encoding step
    x = rans_base_encode_step(x_shrunk,s)
    assert x in [L,H]
    return x, out_bits
```

- Before encoding a symbol into the state `x <- base_step(x_prev,s)` , we shrink the state

  `x_prev -> x_shrunk` , so that `base_step(x_shrunk,s) in [L,H]`

# Streaming rANS: shrink state function

```
def shrink_state(x_prev,s):
    ...
    return x_shrunk, out_bits
```

```
# input state
x_prev = 22 = 10110b

# stream out bits from x_prev until we are sure base_step(x_shrunk,s) in [L,H]
x_shrunk = 10110b = 22, out_bits = ""
x_shrunk = 1011b = 11, out_bits = "0"
x_shrunk = 101b = 5, out_bits = "10"
...
```

# Streaming rANS: shrink state function

```python
def shrink_state(x_prev,s):
    ...
    while rans_base_encode_step(x_shrunk,s) not in Interval[L,H]:
        out_bits.prepend(x_shrunk%2)
        x_shrunk = x_shrunk//2
    x_shrunk = x
    return x_shrunk, out_bits
```

```python
# input state
x_prev = 21 = 10110b

# stream out bits from x_prev until we are sure base_step(x_shrunk,s) in [L,H]
x_shrunk = 10110b = 22, out_bits = ""
x_shrunk = 1011b = 11, out_bits = "0"
x_shrunk = 101b = 5, out_bits = "10"
...
```

# How many bits to stream out?

- We need to choose `[L,H]` so that it is always possible to stream out bits from the state `x_prev` and guarantee that a at some point `base_encode(x_shrunk, s)` `lies in [L,H]`

- For unique decoding, we need to ensure that forthere only one such `x_shrunk` for which `base_encode(x_shrunk, s) lies in [L,H]`

- Condition satisfied for `[L,H] = [M*t, 2*M*t - 1]`

# Streaming rANS: Full Encoding

```python
## streaming params
t = 1 # can be any uint
L = M*t
H = 2*M*t - 1

######################### Streaming rANS Encoder #########################
def rans_base_encode_step(x,s):
    x_next = (x//freq[s])*M + cumul[s] + x%freq[s]
    return x_next

def shrink_state(x,s):
    # initialize the output bitarray
    out_bits = BitArray()

    # shrink state until we are sure the encoded state will lie in the correct interval
    while rans_base_encode_step(x,s) not in Interval[freq[s]*t,2*freq[s]*t - 1]:
        out_bits.prepend(x%2)
        x = x//2
    x_shrunk = x
    return x_shrunk, out_bits

def encode_step(x,s):
    # shrink state x before calling base encode
    x_shrunk, out_bits = shrink_state(x, s)

    # perform the base encoding step
    x_next = rans_base_encode_step(x_shrunk,s)

    return x_next, out_bits
```

# Streaming rANS: Full decoder

```python
########################## Streaming rANS Decoder
def rans_base_decode_step(x):
    ...
    return (s,x_prev)

def expand_state(x_shrunk, enc_bitarray):
    # init
    num_bits_step = 0

    # read in bits to expand x_shrunk -> x
    x = x_shrunk
    while x not in Interval[L,H]:
        x = x*2 + enc_bitarray[num_bits_step]
        num_bits_step += 1
    return x, num_bits_step

def decode_step(x, enc_bitarray):
    # decode s, retrieve prev state
    s, x_shrunk = rans_base_decode_step(x)

    # expand back x_shrunk to lie in Interval[L,H]
    x_prev, num_bits_step = expand_state(x_shrunk, enc_bitarray)
    return s, x_prev, num_bits_step

def decode(encoded_bitarray, num_symbols):
    ...
    return reverse(symbols) # need to reverse to get original sequence
```

# Streaming rANS in practice:

- `M = power of 2`

- `[L,H] = Mt, 2Mt-1` -> `t` chosen as large as possible.

- typically byte/words are streamed out instead of a bit

# Byte/word streaming:

typically byte/words are streamed out instead of a bit

(numbers from https://github.com/rygorous/ryg_rans)

```
# reading/writing one byte at a time
rANS encode: 9277880 clocks, 12.1 clocks/symbol (299.6MiB/s)
rANS decode 14754012 clocks, 19.2 clocks/symbol (188.4MiB/s)

# reading/writing 32 bits at a time:
rANS encode: 7726256 clocks, 10.1 clocks/symbol (359.8MiB/s)
rANS decode: 12159778 clocks, 15.8 clocks/symbol (228.6MiB/s)
```

# Caching rANS computations

- `x in [L,H]` -> for small values interval size, can we cache the `encoding` ?

```python
def encode_step(x,s):
    # shrink state x before calling base encode
    x_shrunk, out_bits = shrink_state(x, s)

    # perform the base encoding step
    x_next = rans_base_encode_step(x_shrunk,s)

    return x_next, out_bits
```

# Caching rANS computation:

- `rans_base_encode_step(x_shrunk,s)` goes from

$$[L, H] \times \{alphabet\} \rightarrow [L, H]$$

**STREAMING-RANS AS A FSE**

Symbol Counts, $\mathcal{F}$ `3,3,2`

Try it

Output State

| Input State | A | B | C |
| --- | --- | --- | --- |
| 8 | 9 | 12 | 14 |
| 9 | 9 | 12 | 14 |
| 10 | 10 | 13 | 14 |
| 11 | 10 | 13 | 14 |
| 12 | 8 | 11 | 15 |
| 13 | 8 | 11 | 15 |
| 14 | 8 | 11 | 15 |
| 15 | 8 | 11 | 15 |

# Caching rANS computation -> tANS

- Both encoding, decoding can be appropriately modified to be basically in terms of lookup-tables

- Lookup-tables are fast, so the encoding/decoding speeds for `tANS` are similar to Huffman coding

- Cannot use a very large interval `[L,H]`, so there is some compression loss (based on how much memory you can use)

# Asymmetric Numeral Systems:

- ANS class of algorithms -> very flexible and can tradeoff compression ratio a bit for very fast implementation:
  ( `rANS` -> `tANS` )

| Codec | Encode speed | Decode speed | compression |
|---|---|---|---|
| Huffman coding | 252 Mb/s | 300 Mb/s | 1.66 |
| Arithmetic coding | 120 Mb/s | 69 Mb/s | 1.24 |
| rANS | 76 Mb/s | 140 Mb/s | 1.24 |
| tANS | 163 Mb/s | 284 Mb/s | 1.25 |

NOTE -> Speed numbers from: Charles Bloom's blog

# rANS, tANS implementations in SCL

- rANS:

  https://github.com/kedartatwawadi/stanford_compression_library/blob/main/scl/compressors/rANS.py

- tANS:

  https://github.com/kedartatwawadi/stanford_compression_library/blob/main/scl/compressors/tANS.py

**Next Class -> Beyond IID distributions, how to deal with correlated sources?**