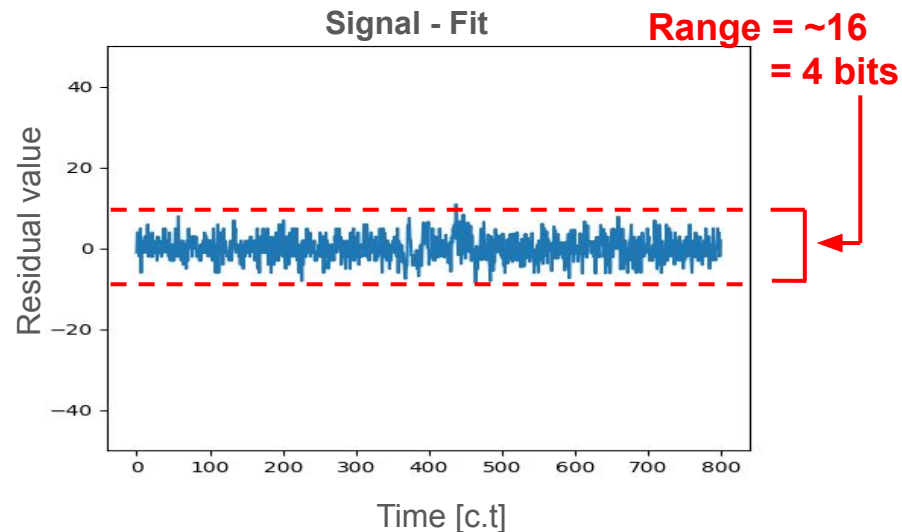
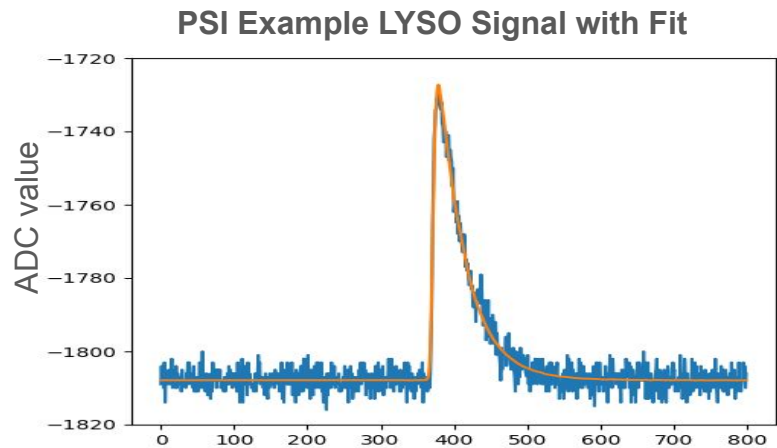


Asymmetric Numeral Systems (ANS) for Compression of Template Fit Data

Jack Carlton
University of Kentucky

Template Fitting

- Can construct a continuous template for our traces $T(t)$
- Can fit traces using template:
$$f(t) = A \cdot T(t - t_0) + B$$
- Storing unfit traces takes ~ 12 bits per ADC sample
- Storing residuals takes ~ 4 bits per ADC sample
- By fitting and compressing, we can expect compress the data by a **factor of ~ 3** for “perfect” fits
 - Effective compressing random noise



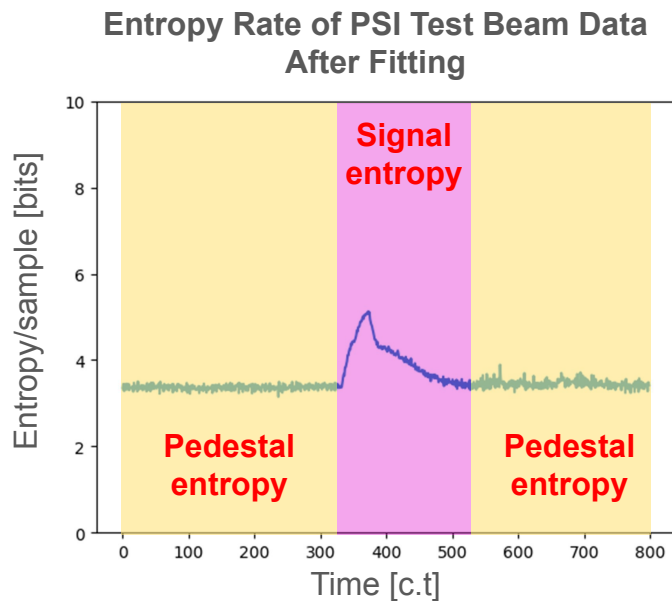
Theoretical Best Compression

- For lossless compression, the best possible compression rate is the entropy rate
- Entropy rate of pedestal part of signal is **3.4 bits per ADC sample**
 - A perfect fit would reduce signal to pedestal noise
- Thus with perfect fit + compression true perfect compression factor is **~3.5**
 - Highly dependent on size of noise

Entropy Rate Formula
(First order, exact for independent residuals)

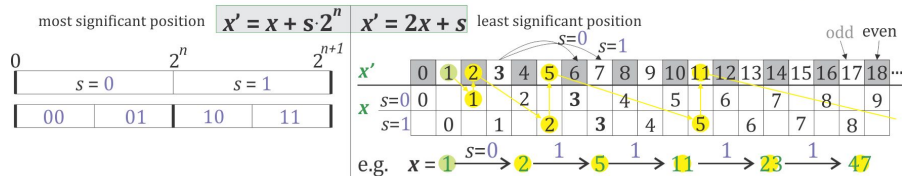
$$H(X_i) = \sum_{\text{traces}} p(X_i) \log_2 (p(X_i))$$

$X_i \equiv$ Random variable for i^{th} ADC sample

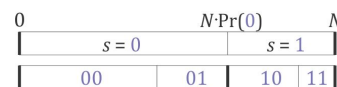


What is Asymmetric Numeral Systems (ANS)?

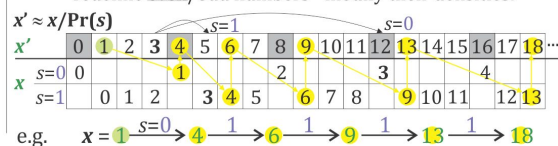
We have information stored in a number x and want to add information of symbol $s=0,1$:
asymmetrize ordinary/symmetric **binary system**: optimal for $\Pr(0)=\Pr(1)=1/2$



range/arithmetic coding:
rescale ranges



some **asymmetric binary system** for $\Pr(0) = 1/4$, $\Pr(1) = 3/4$
 redefine even/odd numbers - modify their densities:



- Entropy encoding compression method
- Manages a state that represents encoded data

- Can separate into many states for GPU parallelization
- Better at encoding “fractional bits” than other methods like Huffman or Rice-Golomb

- [Wikipedia Link](#)

Alphabet and probabilities

$$\mathcal{A} = \{s_1, \dots, s_k\}, \quad p_i = \Pr(s_i), \quad \sum_{i=1}^k p_i = 1$$

Finite precision model

$$F = 2^m, \quad f_i \in \mathbb{N}, \quad \sum_{i=1}^k f_i = F, \quad \frac{f_i}{F} \approx p_i$$

Cumulative frequencies

$$C_i = \sum_{j < i} f_j$$

ANS state

$$x \in \mathbb{N}$$

Encoding update (rANS)

$$x' = \left\lfloor \frac{x}{f_s} \right\rfloor F + (x \bmod f_s) + C_s$$

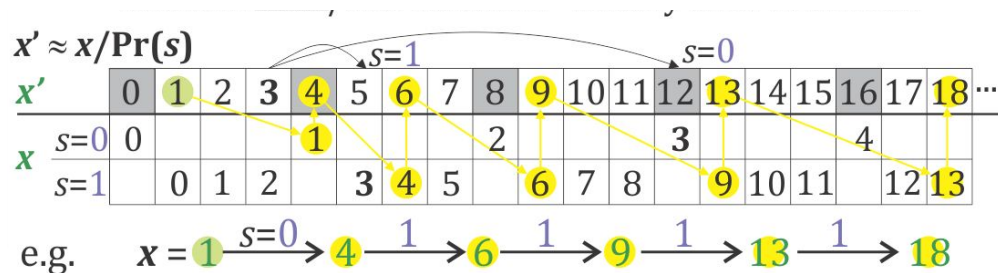
Decoding

$$r = x' \bmod F, \quad s \text{ such that } C_s \leq r < C_s + f_s,$$

$$x = f_s \left\lfloor \frac{x'}{F} \right\rfloor + (r - C_s)$$

What is Table Asymmetric Numeral Systems (tANS)?

- You can define the behavior of ANS (for fixed parameter choices) in a lookup table
- This removes the need for operations, rather each stage is a table lookup
 - Division expensive on GPU
- The cost is needing to store a table
 - This can be built one time based on measured probabilities of samples



Why ANS is good at being parallelized

- Two example approaches:
 - Block parallelism
 - Interleaved parallelism
- Both make it so:
 - One thread holds one state
 - Each state has an identical update rule
 - No synchronization is required across threads
- These make parallelization much more natural than other methods
 - Ex. Rice-Golomb, Huffman

Sequential structure of a single stream

$x_{i+1} = T(x_i, s_i) \Rightarrow$ single dependency chain

Parallelization principle

Maintain K independent states $x^{(0)}, \dots, x^{(K-1)}$

Block-level parallelism

Partition data into blocks B_0, \dots, B_{K-1}

Encode each block with its own ANS state

Interleaved (round-robin) parallelism

For symbol index $i = 0, 1, 2, \dots$: $j = i \bmod K$

$$x^{(j)} \leftarrow T(x^{(j)}, s_i)$$

Explicitly:

$$\begin{aligned} s_0 &\rightarrow x^{(0)}, s_1 \rightarrow x^{(1)}, \dots, s_{K-1} \rightarrow x^{(K-1)}, \\ s_K &\rightarrow x^{(0)}, s_{K+1} \rightarrow x^{(1)}, \dots, s_{2K-1} \rightarrow x^{(K-1)}, \\ s_{2K} &\rightarrow x^{(0)}, s_{2K+1} \rightarrow x^{(1)}, \dots \end{aligned}$$

Things to test/research questions

- At some point Sean had a template fitter for 2023 PSI LYSO waveforms
 - Can we easily get a sample of LYSO data to template fit to produce compressible integer residuals?
- Assuming we have integer residuals, these are the things we want to compress with ANS
 - What data throughput (in GB/s) can we expect using tANS?
 - What compression factor can we expect using tANS?
 - How do the above depend on choice of parameters (ex. Total frequency size F , which is related to the size of the table)
 - Expectation: as F goes up, compression factor improves but performance suffers
 - Compare compression factor to [Shannon Entropy](#), which is the theoretical limit Sean computed for 2023 PSI LYSO data
 - Other useful metrics?


Auxiliary Slides

Should we delta encode?

- **Probably not**
- Delta encoding hurts compression if data has small covariance
- As a result, you should only delta encode if there is a positive correlation

Delta encoded
expectation

non-delta encoded
expectation


$$\mathbb{E}[(r_i - r_{i-1})^2] < \mathbb{E}[r_i^2]$$

$$\mathbb{E}[(r_i - r_{i-1})^2] = \mathbb{E}[r_i^2] + \mathbb{E}[r_{i-1}^2] - 2\mathbb{E}[r_i r_{i-1}]$$

$$\mathbb{E}[(r_i - r_{i-1})^2] < \mathbb{E}[r_i^2] \iff \text{Cov}(r_i, r_{i-1}) > 0$$

Alternatively, look at the correlation coefficient to see “how much” delta encoding will help or hurt

$$\rho = \frac{\mathbb{E}[r_i r_{i-1}]}{\sqrt{\mathbb{E}[r_i^2] \mathbb{E}[r_{i-1}^2]}} > 0$$

Template Fitting - How to store residuals

Model prediction from stored fit parameters

$$\hat{y}(t; \theta) = AT(t - t_0) + B, \quad \theta = (A, t_0, B)$$

True data (integer ADC samples)

$$y(t) \in \mathbb{Z}$$

Define integer prediction via deterministic rounding

$$\tilde{y}(t) = \lceil \hat{y}(t; \theta) \rceil$$

Define integer residual

$$r(t) = y(t) - \tilde{y}(t)$$

Lossless reconstruction

$$y(t) = \tilde{y}(t) + r(t)$$

Equivalently, using the floor identity

$$r(t) = \lfloor y(t) - \hat{y}(t; \theta) \rfloor \in \mathbb{Z}$$

Floor/Ceiling identity (for $y \in \mathbb{Z}$, $x \in \mathbb{R}$):

$$\lfloor y - x \rfloor = y - \lceil x \rceil$$

Therefore

$$y(t) = \lceil \hat{y}(t; \theta) \rceil + \lfloor y(t) - \hat{y}(t; \theta) \rfloor = \lceil \hat{y}(t; \theta) \rceil + r(t)$$

which is exact and lossless.

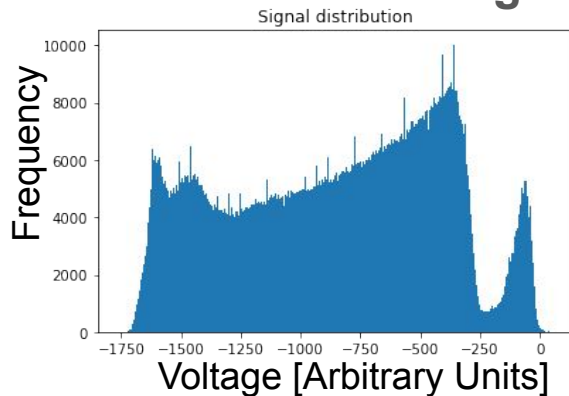
Fit's job:

1. Compute fit params + residuals
2. Store fit params
3. send **integer** residuals to compression algorithm

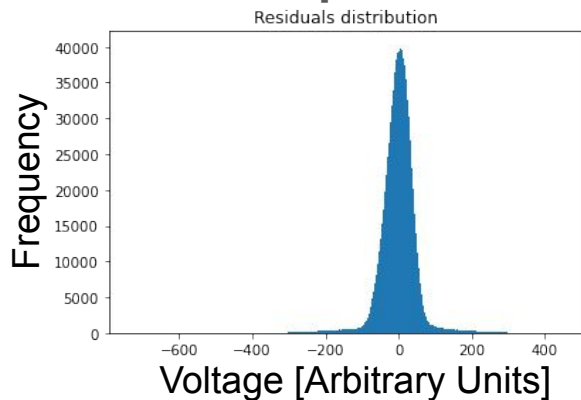
Rounding avoids messy binary float representations of the data

Rounding residuals is still lossless!

No Conditioning



Shape Fit



Higher Order Entropy Estimations

- Assume we have N characters (traces) in our alphabet (data set)

- Zero order:** each character in alphabet is statistically independent $H = \log_2(N)$

- First order:** each character in alphabet is statistically independent, p_i is the probability of that character to occur $H = - \sum_{i=1}^N p_i \log_2(p_i)$

- Second order:** $P_{j|i}$ is correlation between subsequent characters $H = - \sum_{i=1}^N p_i \sum_{j=1}^N P_{j|i} \log_2(P_{j|i})$

- General Model (impractical):** B_n represents the first n characters $H = \lim_{n \rightarrow \infty} \left[-\frac{1}{n} \sum p(B_n) \log_2(B_n) \right]$

What is Asymmetric Numeral Systems (ANS)?

Alphabet and probabilities

$$\mathcal{A} = \{s_1, \dots, s_k\}, \quad p_i = \Pr(s_i), \quad \sum_{i=1}^k p_i = 1$$

Finite precision model

$$F = 2^m, \quad f_i \in \mathbb{N}, \quad \sum_{i=1}^k f_i = F, \quad \frac{f_i}{F} \approx p_i$$

Cumulative frequencies

$$C_i = \sum_{j < i} f_j$$

Normalized ANS state

$$L = 2^R \text{ (normalization bound)}, \quad x \in \mathbb{N}, \quad x \geq L$$

Renormalization (emit/consume bits, base $b = 2^B$)

encode: while $x \geq f_s L$, output $(x \bmod b)$, $x \leftarrow \lfloor x/b \rfloor$

decode: while $x < L$, $x \leftarrow bx + \text{next bits}$

Encoding update (rANS)

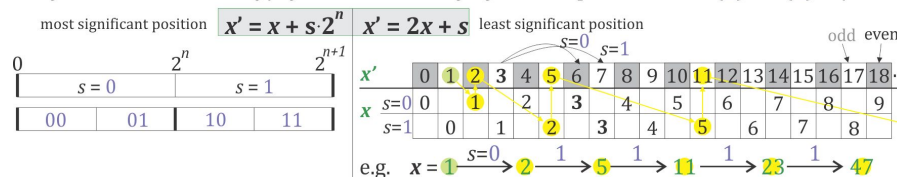
$$x' = \left\lfloor \frac{x}{f_s} \right\rfloor F + (x \bmod f_s) + C_s$$

Decoding

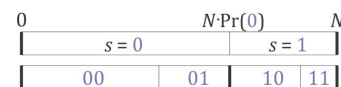
$$r = x' \bmod F, \quad s \text{ such that } C_s \leq r < C_s + f_s,$$

$$x = f_s \left\lfloor \frac{x'}{F} \right\rfloor + (r - C_s)$$

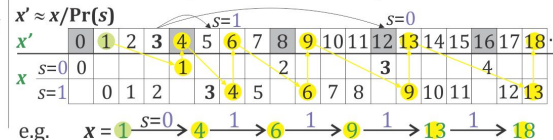
We have information stored in a number x and want to add information of symbol $s=0,1$:
asymmetrize ordinary/symmetric **binary system**: optimal for $\Pr(0)=\Pr(1)=1/2$



range/arithmetic coding:
 rescale ranges



some **asymmetric binary system** for $\Pr(0) = 1/4, \Pr(1) = 3/4$
 redefine even/odd numbers - modify their densities:



This is like a “buffer” between bitstream and the state. Your state, likely encoded in a 64 bit integer, will overflow when encoding or decoding. So you need to “emit” or “consume” bits to an arbitrary long bitstream as that happens.