

---

# Serializing Nouns

N. E. Davis ~lagrev-nocfep, Brian Klatt ~zod,\*  
& Sam Parker ~zod†  
Zorp Corp, Zorp Corp, & –

## Abstract

Noun serialization is commonly used for Nock communication, both between instances like Urbit ships and with the runtime and the Unix host operating system. This article describes and compares the two principal conventions for representing nouns in slightly compressed form as byte arrays, as well as introduces variant encodings for educational purposes.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Naïve Serialization</b>	<b>3</b>
<b>3</b>	<b>Practical Serialization</b>	<b>10</b>
3.1	newt Encoding . . . . .	14
<b>4</b>	<b>Directed Graph Encoding</b>	<b>14</b>

---

\*Brian Klatt contributed the description of +jam.

†Sam Parker contributed the description of +bulk.

We gratefully acknowledge the contributions of the original architects of +jam/+cue as well.

<b>5</b>	<b>Aligned Serialization</b>	<b>14</b>
<b>6</b>	<b>Benchmarks</b>	<b>15</b>
<b>7</b>	<b>Conclusion</b>	<b>15</b>

## 1 Introduction

The Nock combinator calculus deals in nouns and does not know about bit encodings or memory layouts. A Nock interpreter (runtime) must, however, deal with the practicalities. In other words, there must be a way of writing down an abstract binary tree (consisting of cells/pairs and atoms) as an actual, physical array of bits in memory. Every possible Nock noun can be represented as a finite sequence of bytes (an atom), and there are multiple ways to do so.<sup>1</sup>

A noun serialization strategy is rather like a Gödel numbering in that it systematically encodes a mathematical object (a noun) as a number (an atom). Unlike Gödel numbering, which classically serially encodes the symbols of mathematical statements, noun serialization encodes a binary tree structure.<sup>2</sup> Because a noun may be any atom—and atoms cannot have leading zeroes—both structure and value need to be unambiguously encoded and cannot be simply delimited (as by a 0 bit or similar). There are two basic strategies to encode a noun as an atom:

1. Run-length serialization, with or without references.
2. Directed graph serialization, depending on a reentrant graph encoding.

Both of these embody a tradeoff between simplicity, size, and speed. Below, we describe both of these strategies and some new variants which offer possible advantages. Reversible noun encoding is essential for Nock communication, both between Nock execution layers such as Vere and NockVM, and between the runtime and the host operating system.

---

<sup>1</sup>The converse is not true: not every atom represents a valid deserialization or conversion from a graph encoding.

<sup>2</sup>This also echoes the way that S-expressions are encoded in Lisp.

## 2 Naïve Serialization

A noun may be an atom or a cell, and so we first propose to label nodes by type, with 0 for atom and 1 for cell.

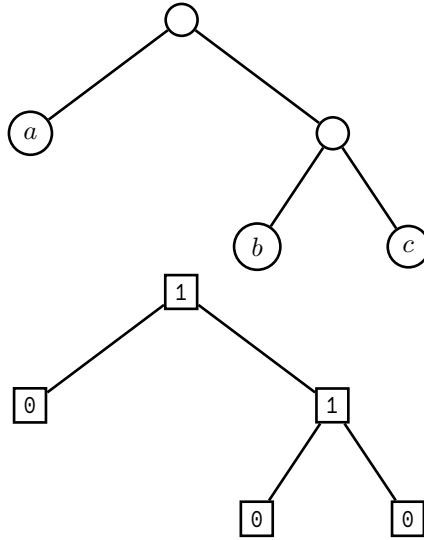


Figure 1: The noun [a b c] (the same as [a [b c]]) as a binary tree and its metadata labeling in the simplest scheme.

The noun in Figure 1 is represented as a binary tree, where each node is labeled with its type (atom or cell) and its value (if applicable). To serialize this noun, we can perform a pre-order traversal of the tree, recording the type and value of each node. The serialized form would be (in LSB order):

---

c 0 b 0 1 a 0 1

---

The issue, however, is that in decoding this bitstream we do not know where an atom value (such as a) ends. After all, the point of serialization is to be able to deserialize the bitstream back into the original noun. Thus we need a way to indicate the length of each atom.

Let’s augment this scheme by providing length information for each atom. One basic idea is to include a string of 0 bits as the least-significant bits (LSB) of the atom which indicates the length of the atom in bits. Since the atom could be zero, we also delimit the length with a 1 bit. (See Table 1 for examples.)

$$a = a_{\ell-1}a_{\ell-2} \dots a_1a_0 \xrightarrow[\text{encode}]{} \bar{a} = a_{\ell-1}a_{\ell-2} \dots a_1a_0 \ 1 \ 0^n,$$

where  $\bar{a}$  is the encoded atom  $a$  and  $a$  has length  $\ell$  bits.<sup>3</sup> This is a unary run-length encoding (RLE) using zeros.<sup>4</sup>

Table 1: Example atom encodings with length information appended as trailing zeros (unnamed encoding).

Atom $a$	Length (bits)	Encoded Form $\bar{a}$
0b0	1	0b1
0b1	1	0b110
0b10	2	0b10100
0b11	2	0b11100
0b101	3	0b1011000
0b111	3	0b1111000
0x70	7	0b111000010000000
0xff	8	0b11111111100000000

We can make one further optimization: since the most significant bit (MSB) of the atom is always 1 except in the case of zero (which is 0b10 with the leading zero trimmed), we can omit it from the encoded form. This saves one bit for every non-zero atom. Let’s call this scheme *+fat* encoding. (See Table 2 for examples.)

#### *+fat* Serialization Algorithm

- **Encode.** Write<sup>5</sup>  $\ell$  bits of 0, where  $\ell$  is the length of  $a$  in

---

<sup>3</sup>We will use the overbar notation  $\bar{a}$  throughout to notionally indicate the encoded form of  $a$ , regardless of the method used.

<sup>4</sup>Zero is something of a “special case”: it has no length and leading zeros are trimmed, so its encoding consists only of the delimiter, 0b1.

<sup>5</sup>We always say this from the LSB, so from the right as written on the page.

unary, followed by 1. The next  $\ell$  bits are the atom  $a$  in binary without the leading 1.

- **Decode.** Read the unary length  $\ell$  by counting 0 bits until the first 1 bit. Read the next  $\ell$  bits to get the atom  $a$  without the leading 1. Prepend a 1 bit to  $a$  unless  $\ell$  is zero, in which case the atom is zero.

Table 2: Example atom encodings (+fat encoding).

Atom $a$	Length (bits)	Encoded Form $\bar{a}$
0b0	1	0b1
0b1	1	0b10
0b10	2	0b100
0b11	2	0b1100
0b101	3	0b11000
0b111	3	0b111000
0x70	7	0b11000010000000
0xff	8	0b1111111100000000

The primary problem with the +fat scheme is that the encoding of an atom is twice the bit length of the original atom. The next optimization is to use the bit length of the bit length, which is logarithmic in the bit length.

$$\begin{aligned}
 a &= \underbrace{a_{\ell-1}a_{\ell-2} \dots a_1a_0}_{\ell = \underbrace{\ell_{h-1}\ell_{h-2} \dots \ell_1\ell_0}_h} \\
 a &= a_{\ell-1}a_{\ell-2} \dots a_1a_0 \xrightarrow{\text{+fat encode}} \\
 \bar{a} &= \underbrace{a_{\ell-1}a_{\ell-2} \dots a_1a_0}_{\ell} \underbrace{\ell_{h-2}\ell_{h-1} \dots \ell_1\ell_0}_{h-1} \underbrace{1\ 0 \dots 0}_h,
 \end{aligned}$$

where  $\ell$  is the bit length of  $a$  and  $h$  is the “hyper bit length” of  $a$ . This yields a bit length encoding of  $1 + 2\log_2(\ell)$  bits rather than  $2\ell$  bits. This scheme is called +mat encoding, and it is currently in use as an internal auxiliary format for other serializations.

#### `+mat` Serialization Algorithm

- **Encode.** Write the number  $h$  of 0 bits in unary, followed a leading 1 and then  $h - 1$  bits for the low bits of  $\ell$ . The next  $\ell$  bits are the atom  $a$  in binary.
- **Decode.** Read the unary length of the length  $h$  by counting 0 bits (from the LSB/right) until the first 1 bit. Read the next  $h - 1$  bits to get  $\ell - 1$ , the bit length of the atom. Finally, read the next  $\ell$  bits (adding the leading 1) to get the atom  $a$ .

Some worked examples of `+mat` encoding are included in Table 3.

Table 3: Example atom encodings (+mat encoding). The encoded form is the same result as  $q: (\text{mat } a)$  in Hoon.

Atom $a$	Length (bits)	Mathematics	Encoded Form $\bar{a}$
0b0	$\ell = 1$	$\underbrace{a_{\ell-1} \dots a_0}_{\ell=0} \underbrace{1 \ 0 \dots 0}_{h=0}$	0b1
0b1	$\ell = 1$	$\underbrace{1 \ \ell_{h-2} \dots \ell_0}_{\ell=1} \underbrace{1 \ 0}_{h=1}$	0b110
0b10	$\ell = 2$	$\underbrace{10}_{a=2 \text{ low bits of } \ell} \underbrace{0}_{h=2} \underbrace{1 \ 00}_{h=2}$	0b100100
0b11	$\ell = 2$	$\underbrace{11}_{a=3 \text{ low bits of } \ell} \underbrace{0}_{h=2} \underbrace{1 \ 00}_{h=2}$	0b110100
0b100	$\ell = 3$	$\underbrace{100}_{a=4 \text{ low bits of } \ell} \underbrace{1}_{h=2} \underbrace{1 \ 00}_{h=2}$	0b1001100
0b101	$\ell = 3$	$\underbrace{101}_{a=5 \text{ low bits of } \ell} \underbrace{1}_{h=2} \underbrace{1 \ 00}_{h=2}$	0b1011100
0b1111	$\ell = 4$	$\underbrace{1111}_{a=15 \text{ low bits of } \ell} \underbrace{00}_{h=3} \underbrace{1 \ 000}_{h=3}$	0b1111001000
0x70	$\ell = 7$	$\underbrace{1110000}_{a=112 \text{ low bits of } \ell} \underbrace{11}_{h=3} \underbrace{1 \ 000}_{h=3}$	0b1110000111000

With `+mat` encoding in hand, we can produce a proper serialization scheme for nouns by combining the type bit with the encoded atom or recursively encoding cells. We call this scheme `+gel` encoding.

As before, we use 0 for atom and 1 for cell. If the noun is an atom, we follow the type bit with the `+mat` encoding of the atom. If the noun is a cell, we follow the type bit with the `+gel` encoding of the head and then the `+gel` encoding of the tail.

Consider the noun `[a [b c]]`, represented as a binary tree in Figure 2. Reading across in a depth-first traversal, we produce an encoding order:

```
(mat 1)(0)(mat 0)(0)(1)(mat 4)(0)(1)
```

We already have the `+mat` encodings from Table 3; substituting these in yields:

```
[110](0)[1](0)(1)[1001100](0)(1)
0b110.0.1.0.1.1001100.0.1
0b1100101100110001
```

### `+gel` Serialization Algorithm

- **Encode.**

- If the noun is an atom, write 0 followed by the `+mat` encoding of the atom.
- If the noun is a cell, write 1 followed by the `+gel` encoding of the head and then the `+gel` encoding of the tail.

- **Decode.** Read the first bit.

- If it is 0, read the next bits as a `+mat` encoded atom.
- If it is 1, recursively read the next bits as the head and then the tail.

`+gel` encoding is straightforward to read and write even by hand. Worked examples are included in Table 4.

TODO opposite of gel as goo

Our code implementation for `+gel` is as follows:



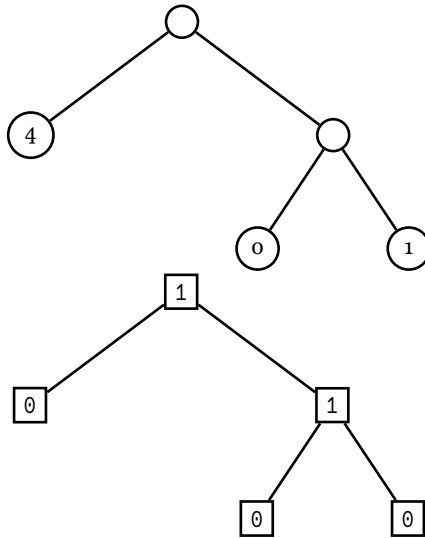


Figure 2: The noun `[4 0 1]` as a binary tree and its metadata labeling in the simplest scheme. Compare Figure 1.

---

```

! : | %
++ gel
  =gel ! : | = a = *
    ^ - @
5  =+ l=0
    =+ b=0
    =< -
    | -
    ?^ a
10  =+ lv=$(a -.a)
    =+ rv=$(a +.a)
    =+ [c l]=(mash rv lv)
      [(con (lsh [0 1] c) 0b1) +(1)]
    ? : =(0 a) [0b10 4]
15  :: need another mash in here for unary length
    =+ [c l]=(mash [a l] [b +((met 0 b))])
      [(con (lsh [0 1] c) 0b0) +(1)]
    :: length of atom in unary

```

Table 4: Example noun encodings (+gel encoding). . dot delimits bit fields for readability.

Noun	Encoded Form
[0 0]	0b1.0.1.0.1
[1 0]	0b110.0.1.0.1
[2 1 0]	0b110.0.1.0.1.100100.0.1
[[2 3] 1 0]	0b110.0.1.0.1.110100.0.100100.0.1.1

```

++ len
20 | = a=@
    ^ - @
    (fil 0 (met 0 a) 0b1)
:: mash two atoms together
++ mash
25 | = [a=[p=@ l=@] b=[p=@ l=@]]
    ^ - [c=@ l=@]
    :- (con (lsh [0 1.b] p.a) p.b)
    (add 1.a 1.b)
--
    
```

---

### 3 Practical Serialization

Whatever the pedagogical advantages of +gel, the algorithm is still relatively verbose. Nock-based operating functions have preferred in practice to use +jam encoding, which is a more compact run-length encoding with subtree deduplication.<sup>6</sup> +jam improves the basic strategy by altering the RLE algorithm slightly and supporting internal references for noun subtrees that have already been encoded.

+jam converts a noun into a buffer and deduplicates repeated subtrees. It walks subtrees and encodes each in a way that allows for efficient storage and retrieval, while also permitting references to previously encoded values. Basically, it

---

<sup>6</sup>Note the calculation of ~dozreg-toplud, p. TODO of this issue, that an operational Arvo instance may have up to  $1.66 \times 10^{21}$  nouns, greatly reduced by structural sharing.

takes `+gel` and adds a “caching mode” flag; instead of repeating subtrees, encode the index in the full jammed noun to where the current node’s jammed subnoun was first written. Thus the code changes from 0 for atom and 1 for cell to:

- 0 for atom.
- 01 for cell.
- 11 for cached node.

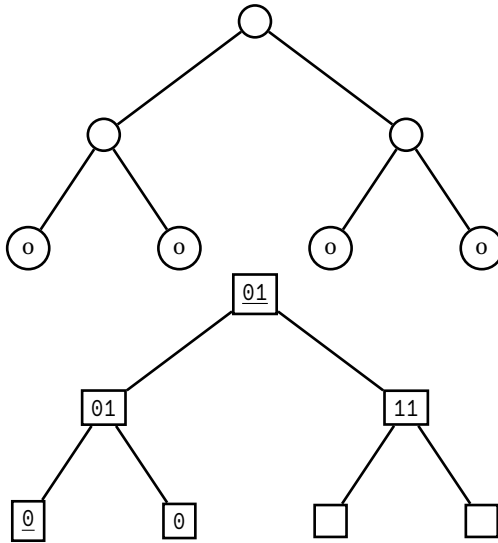


Figure 3: The noun `[[0 0] 0 0]` as a binary tree and its meta-data labeling for the `+jam` scheme.

Following the metadata labels given in Figure 3, the encoding order for `[[0 0] 0 0]` is:

```

(mat 2)(11)(mat 0)(0)(mat 0)(0)(01)(01)
0b100100.11.1.0.1.0.01.01
0b1001001110100101

```

The cached entry (`mat 2`) indicates that the cached subtree is a copy of the subtree located at tree address 2 in the overall jammed noun.

However, this bit of cleverness isn't the end: `+jam` also decides *when* to cache subtrees based on the size of the nouns involved. For instance, `[3 3 3]` encodes as:

```
(mat 3)(0)(mat 3)(0)(01)(mat 3)(0)(01)
0b110100.0.110100.0.01.110100.0.01
0b1101000110100001110100001
```

wherein no caching is used. When 3 is first encountered at index 2, 2 and 3 have the same bit lengths; since `(mat 2)` and `(mat 3)` are the same length, caching would not save any space and would add decoding overhead due to dereferencing.

In contrast, `[4 4 4]` encodes as:

```
(mat 2)(11)(mat 2)(11)(01)(1001100)(0)(01)
0b100100.11.100100.11.01.1001100.0.01
0b1001001110010011011001100001
```

where caching is used for the second and third occurrences of 4. Here, 4 has bit length 3 which is greater than the bit length of 2 which is 2, and caching saves space overall in the jammed noun.

### `+jam` Serialization Algorithm

- **Encode.**

- If the noun is an atom, write 0 followed by the `+mat` encoding of the atom.
- If the noun is a cell, check if it has been previously encoded.
  - \* If it has, write 11 followed by the `+mat` encoding of the index where it was first encoded.
  - \* If it has not, write 01 followed by the `+jam` encoding of the head and then the `+jam` encoding of the tail. Store the current index for future reference.

- **Decode.** Read the first bit.

- If it is 0, read the next bits as a +mat encoded atom.
- If it is 1, read the next bit.
  - \* If it is 0, recursively read the next bits as the head and then the tail.
  - \* If it is 1, read the next bits as a +mat encoded index and retrieve the cached subtree.

Here is the reference implementation of +jam in Hoon:

---

```

++ jam
~/ %jam
|= a=*
^ - @
5  += b=0
   += m=`(map * @)`~
   =< q
   |- ^- [p=@ q=@ r=(map * @)]
   += c=(~(get by m) a)
10  ?~ c
     => .(m (~(put by m) a b))
     ? : ?=(@ a)
         += d=(mat a)
           [(add 1 p.d) (lsh 0 q.d) m]
15  => .(b (add 2 b))
     += d=$(a -.a)
     += e=$(a +.a, b (add b p.d), m r.d)
     :+ (add 2 (add p.d p.e))
       (mix 1 (lsh [0 2] (cat 0 q.d q.e)))
20  r.e
     ? : ?&(?=(@ a) (lte (met 0 a) (met 0 u.c)))
         += d=(mat a)
           [(add 1 p.d) (lsh 0 q.d) m]
         += d=(mat u.c)
25  [(add 2 p.d) (mix 3 (lsh [0 2] q.d)) m]

```

---

A Python implementation of +jam is included in Appendix A.  
Here is the reference implementation of +cue in Hoon:

---

TODO

---

Likewise, a Python implementation of +cue is included in  
Appendix A.

### 3.1 newt Encoding

“Newt” encoding is a runtime-oriented extension of +jam based noun serialization which adds a short identifying header in case of future changes to the runtime’s serialization format. A version number (currently a single bit) precedes a RLE serialization length followed by the +jam serialization of the noun. The version number is currently 0b0.

---

V.LLLL.JJJJ.JJJJ.JJJJ.JJJJ.JJJJ.JJJJ

---

where V is the version number, L is the total length of the noun in bytes, and J is the +jam serialization of the noun.

Runtime communications vanes like %khan and %lick utilize this encoding locally. It is exclusively used as a host OS runtime affordance at the current time.

## 4 Directed Graph Encoding

A directed graph encoding has been independently proposed twice, once by Tlon in the original Hoon codebase as a +\$silo encoding and once by Sam Parker as the +bulk encoding.

## 5 Aligned Serialization

One of the advantages of +jam is its compactness. However, this comes at the cost of speed, since bit-level operations are required to +cue the noun back from its serialized form. If a slightly larger size is acceptable, a byte-aligned serialization could facilitate certain kinds of external inspection without requiring deserialization. (For instance, a byte-aligned head tag could be read for a rapid decision without needing to +cue the entire noun.)

We propose a strategy to modify +jam to align to bytes by padding the length of entries to the nearest byte boundary and marking the distance with a clever binary scheme rather than simply unary. This approach, called +honey, aims to balance compactness and speed for certain use cases while retaining a large degree of conceptual backwards compatibility. (The

change in byte alignment of course breaks strict compatibility.)

There are two fundamental issues for byte alignment: atoms and lengths. Atoms can be padded with leading zeros to the nearest byte boundary without changing their value. Lengths, however, require a new encoding scheme to compensate for the adjustment in expected bit widths.

## 6 Benchmarks

## 7 Conclusion

### Appendix A: Python +jam/+cue

The following is a simple Python implementation of +jam serialization drawn from Urbit's auxiliary pynoun library.

```
from bitstring import BitArray
noun = int | Cell
# The Cell class represents an ordered pair of two nouns.

5 def jam_to_stream(n: noun, out: BitArray):
    """jam but put the bits into a stream

    >>> s = BitArray()
    >>> jam_to_stream(Cell(0,0), s)
10 >>> s
    BitArray('0b100101')
    """

    cur = 0
15 refs = {}

    def bit(b: bool):
        nonlocal cur
        out.append([b])
        cur += 1
20

    def zero():
        bit(False)
```

```

25     def one():
        bit(True)

    def bits(num: int, count: int):
        nonlocal cur
30        for i in range(0, count):
            out.append([(num & (1 << i)) != 0])
            cur += count

    def save(a: noun):
35        refs[a] = cur

    def mat(i: int):
        if 0 == i:
            one()
40        else:
            a = i.bit_length()
            b = a.bit_length()
            above = b + 1
            below = b - 1
45            bits(1 << b, above)
            bits(a & ((1 << below) - 1), below)
            bits(i, a)

    def back(ref: int):
50        one()
        one()
        mat(ref)

    def r(a: noun):
55        dupe = refs.get(a)
        if deep(a):
            if dupe:
                back(dupe)
            else:
60                save(a)
                one()
                zero()
                r(a.head)
                r(a.tail)
65        elif dupe:

```



```
        isize = a.bit_length()
        dsize = dupe.bit_length()
        if isize < dsize:
            zero()
70         mat(a)
        else:
            back(dupe)
    else:
        save(a)
75         zero()
        mat(a)
    r(n)

def jam(n: noun):
80     """urbit serialization: * -> @

    >>> jam(0)
    2
    >>> jam(Cell(0,0))
85     41
    >>> jam(Cell(Cell(1234567890987654321, \\
    ...             1234567890987654321), \\
    ...             Cell(1234567890987654321, \\
    ...             1234567890987654321)))
90     22840095095806892874257389573
    "" ""

    out = BitArray()
    jam_to_stream(n, out)
95     return read_int(len(out), out)

def cue_from_stream(s: BitArray):
    """cue but read the bits from a stream

100     >>> s = BitArray('0b01')
    >>> cue_from_stream(s)
    0
    "" ""

105     refs = {}
    cur = 0
    position = 0
```

```

def bits(n: int):
110     nonlocal cur, position
        cur += n
        result = 0
        for i in range(n):
            result |= (1 if s[position] else 0) << i
115         position += 1
        return result

def one():
    nonlocal cur, position
120     cur += 1
    bit = s[position]
    position += 1
    return bit

125 def rub():
    z = 0
    while not one():
        z += 1
    if 0 == z:
130         return 0
    below = z - 1
    lbits = bits(below)
    bex = 1 << below
    return bits(bex ^ lbits)

135 def r(start: int):
    ret = None
    if one():
        if one():
140             ret = refs[rub()]
        else:
            hed = r(cur)
            tal = r(cur)
            ret = Cell(hed, tal)
145     else:
        ret = rub()
        refs[start] = ret
        return ret
    return r(cur)

```

```
150 def cue(i: int):  
    """urbit deserialization: @ -> *  
  
    >>> str(cue(22840095095806892874257389573))  
155 '[[1234567890987654321 1234567890987654321]  
    1234567890987654321 1234567890987654321]'  
    "" "  
  
    bits = BitArray()  
160 while i > 0:  
        bits.append([i & 1 == 1])  
        i >>= 1  
    return cue_from_stream(bits)
```