
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

1	Introduction	2
2	Naïve Serialization	3
3	Practical Serialization	5
3.1	newt Encoding	6
4	Directed Graph Encoding	6
5	Aligned Serialization	6

*Brian Klatt contributed the description of `++jam`.

†Sam Parker contributed the description of `++bulk`.

6	Benchmarks	7
7	Conclusion	7

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.¹

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.² 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. In this article, we describe both of these strategies and some new variants which offer possible advantages.

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

²This also echoes the way that S-expressions are encoded in Lisp.

2 Naïve Serialization

The simplest way to encode a noun as a binary tree in an atom (byte array) without compression is to utilize bits to mark atom (0) vs. cell (1), the length of the atom in unary of 1s terminated by 0, and the actual value. (Since the leading zeros are stripped, they may need to be added back in.) This is in least-significant byte (LSB) order, so you should read the written atom starting from the *right* in binary.

Thus the atom 0b0 would simply be encoded as:

```
:: 0x0 = 0 for atom; 1 for length (special case);
::          0 for end of length; 0 for value
> `@ub`(jel 0b0)
0b10
```

interpreted as (rightmost, LSB) 0 for atom, run length of 1, and value 0 (stripped). I.e., 0b010 or

←

① ② ③

Other values include:

```
:: 0x1 = 0 for atom; 1 for length;
::          0 for end of length; 1 for value
> `@ub`(jel 0b1)
0b1010
5
:: [0x0 0x1] = 1 for cell; zero, then one
> `@ub`(jel [0b0@ 0b1])
0b1.010@0.010@1
10
:: [0x1 0x0] = 1 for cell; one, then zero
> `@ub`(jel [0b0@ 0b1])
0b10@1.010@1
15
:: 0x2 = 0 for atom; 2 for length;
::          0 for end of length; 2 for value
> `@ub`(jel 0b10)
0b10.0111
20
:: 0xff
> `@ub`(jel 0xff)
```

0b11.1111.1101.1111.1110

```
25 :: [[0x0 0x1] [0x2 0x3]]
> `@ub`(jel [@@b0@ 0b1])
0b10@1.010@1
```

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

```
! : |%
++ jel
=jel !: |= 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 1] [b +((met 0 b))])
[(con (lsh [0 1] c) 0b0) +(1)]
:: length of atom in unary
++ 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 l.b] p.a) p.b)
(add l.a l.b)
--
```

Note that `l` is not the length of an atom in unary, but the length of the encoded noun in binary.

3 Practical Serialization

Whatever the pedagogical advantages of `++jel`, the algorithm has practical flaws: it is subject to collisions XXX and it is verbose.³ `++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.

- get Bryan notes

special-cases o

The new RLE calculation is to post the number plus one in binary rather than unary (e.g., for the length would be `0b010`).

(Since there is a value less significant than the length, the leading zero is not lost but serves as a divider.)

```

++ jam
~/ %jam
|= a=*
^-
5 =+ b=0
=+ m=`(map * @)`~
=< q
|- ^- [p=@ q=@ r=(map * @)]
=+ c=(~(get by m) a)
?~ 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

```

³Note the claim of `~dozreg-toplud`, p. TODO of this issue, that an operational Arvo instance may have up to 1.66×10^{21} nouns, reduced by structural sharing.

```
? : ?&(?=(@ 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)
[(add 2 p.d) (mix 3 (lsh [0 2] q.d)) m]
```

A Python example of `++jam` 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 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 vane 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)
    >>> s
    BitArray('0b100101')
    """

10
    cur = 0
    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
    for i in range(0, count):
        out.append([(num & (1 << i)) != 0])
    cur += count
30

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

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

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

def r(a: noun):
    dupe = refs.get(a)
    if deep(a):
55        if dupe:
```

```
        back(dupe)
    else:
        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()
            mat(a)
70    else:
        back(dupe)
    else:
        save(a)
        zero()
        mat(a)
75
    r(n)

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

    >>> jam(0)
    2
    >>> jam(Cell(0,0))
    41
    >>> jam(Cell(Cell(1234567890987654321, @@
    ...           1234567890987654321), @@
    ...           Cell(1234567890987654321, @@
    ...           1234567890987654321)))
85    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):
    nonlocal cur, position
    cur += n
    result = 0
    for i in range(n):
        result |= (1 if s[position] else 0) << i
    position += 1
    return result

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

125      def rub():
    z = 0
    while not one():
        z += 1
    if 0 == z:
        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():
            ret = refs[rub()]
    140      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)
```