

Horn Binary Serialization Analysis

HCVS 2016

3rd Workshop
on Horn Clauses
for Verification and Synthesis

Gabriele Paganelli
<https://gapag.noblogs.org/>
gapag@distruzione.org

Length	Type	<i>Data</i>	CRC
0	3 4	7 8	N N+1 N+4

Type	<i>Data</i>	Length	CRC
0	3 4	N N+1	N+4 N+5 N+8



Cat.png

Contribution

- **Given** a layout specification,

Length	Type	<i>Data</i>	CRC
0	3 4	7 8	N N+1 N+4

Is there a parser

that can parse any instance stream?

Or: *is the layout deserializable?*

***In practice:* describe a left-to-right parser behaviour using Horn clauses and forward chaining**

1: Formalize a layout specification

- $f [Id]$: fixed length field
- $v [Id]$: variable length field (*varfield*)
- $(\underset{\text{offset}}{Id'} \rightarrow \underset{\text{span}}{\beta}) [Id]$: pointer field

(Length→4) Length

Length	Type	<i>Data</i>	CRC
0	3 4	7 8	N N+1 N+4
	f	v	f

2: Give a name to all fields

(Length \rightarrow 4) Length₀ f₁ v₂ f₃

Length	Type	<i>Data</i>	CRC
0	3 4	7 8	N N+1 N+4

3: Formalize parser's knowledge

- ***The parser knows...***
- *Beg(i)* : where field *i* begins
- *Len(i)* : field *i*'s length
- *Ptr(o,s,i)*: field *i* is a pointer,
with offset at *o*
and spanning *s* fields
- *Val(i)*: field *i*'s contents.

4a: Formalize parser's behaviour

$(\text{Length} \rightarrow 4) \text{Length}_0 \ f_1 \ v_2 \ f_3$

Length	Type	<i>Data</i>	CRC
0	3 4	7 8	N N+1 N+4

True \Rightarrow *Beg(0)*

True \Rightarrow *Len(0)*

True \Rightarrow *Len(1)*

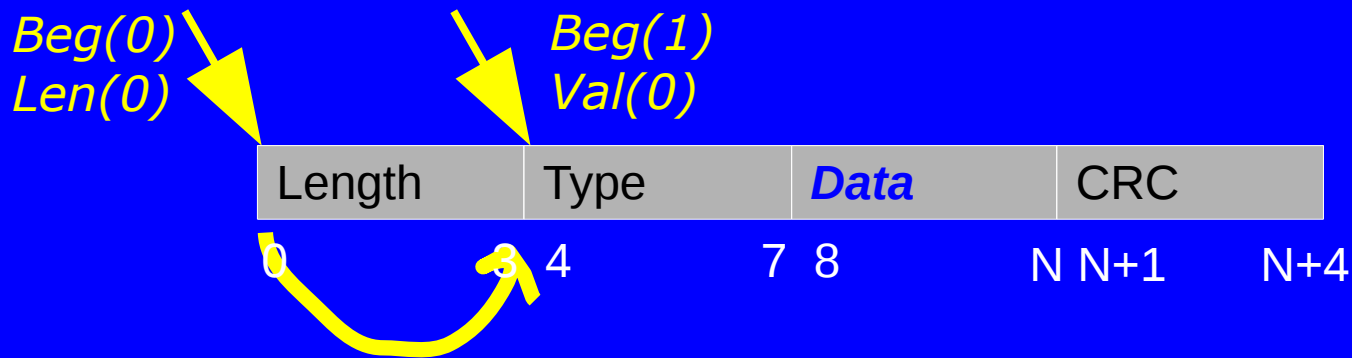
True \Rightarrow *Len(3)*

True \Rightarrow *Ptr(0,4,0)*

4b: Formalize parser's behaviour

- $Beg(i) \wedge Len(i) \Rightarrow Beg(i+1) \wedge Val(i)$ forward
- $Beg(i+1) \wedge Len(i) \Rightarrow Beg(i) \wedge Val(i)$ backward

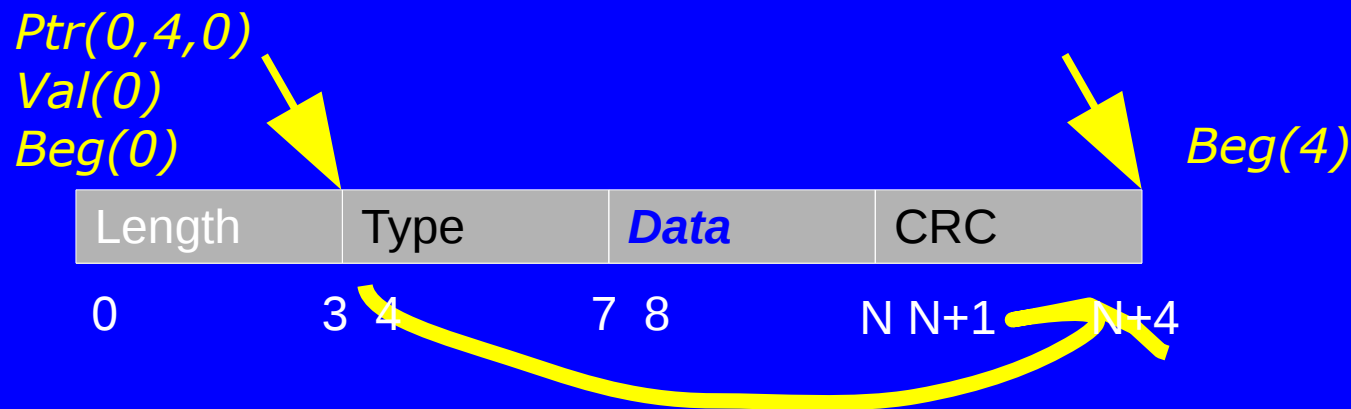
- Read a field backward or forward.



4c: Formalize parser's behaviour

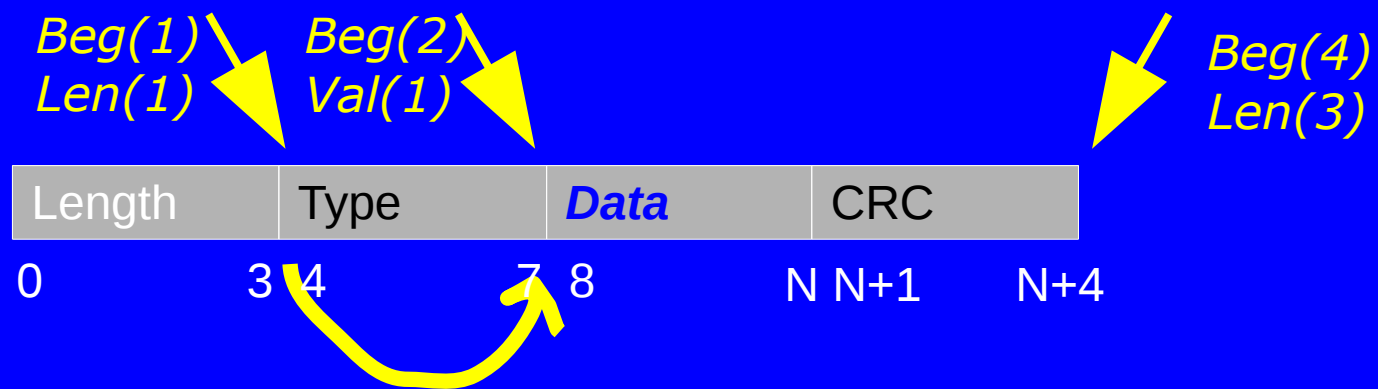
- $Ptr(o,s,i) \wedge Val(i) \wedge Beg(o) \Rightarrow Beg(o+s)$ Jump right
- $Ptr(o,s,i) \wedge Val(i) \wedge Beg(o+s) \Rightarrow Beg(o)$ Jump left

- Follow a pointer backward or forward.



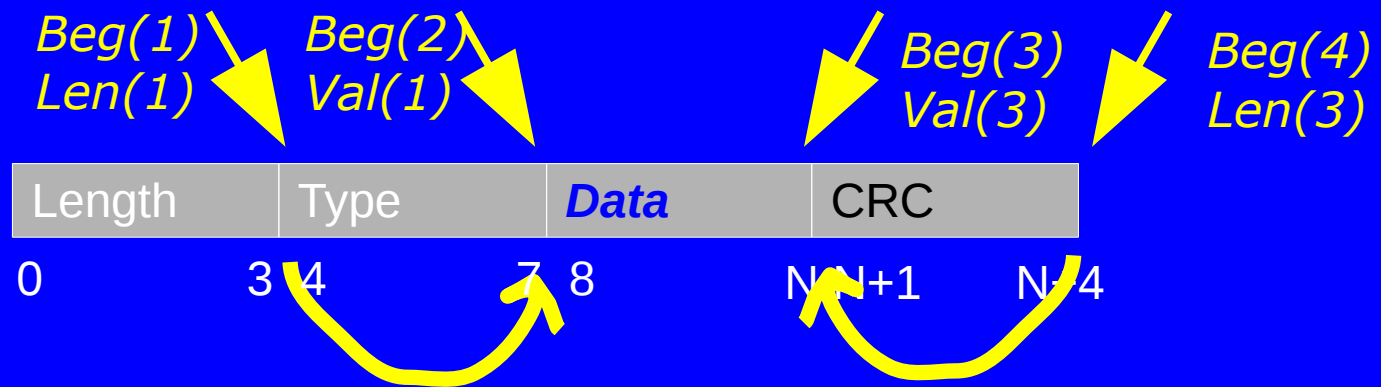
4cc:(example, continued)

- $Beg(i) \wedge Len(i) \Rightarrow Beg(i+1) \wedge Val(i)$ forward
 - $Beg(i+1) \wedge Len(i) \Rightarrow Beg(i) \wedge Val(i)$ backward
-
- Read a field backward or forward.



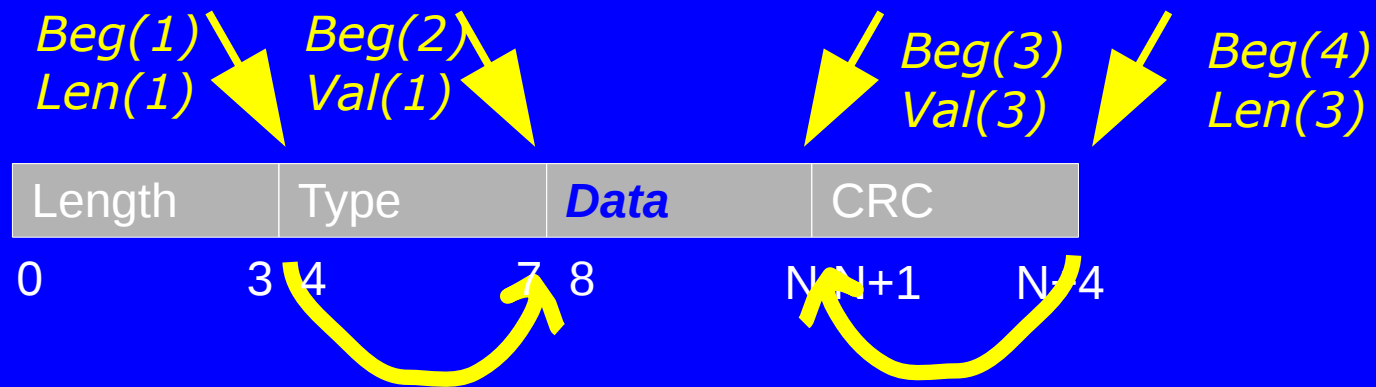
4cc:(example, continued)

- $Beg(i) \wedge Len(i) \Rightarrow Beg(i+1) \wedge Val(i)$ forward
 - $Beg(i+1) \wedge Len(i) \Rightarrow Beg(i) \wedge Val(i)$ backward
-
- Read a field backward or forward.



4cc:(example, continued)

- $Beg(i) \wedge Len(i) \Rightarrow Beg(i+1) \wedge Val(i)$ forward
 - $Beg(i+1) \wedge Len(i) \Rightarrow Beg(i) \wedge Val(i)$ backward
-
- Read a field backward or forward.

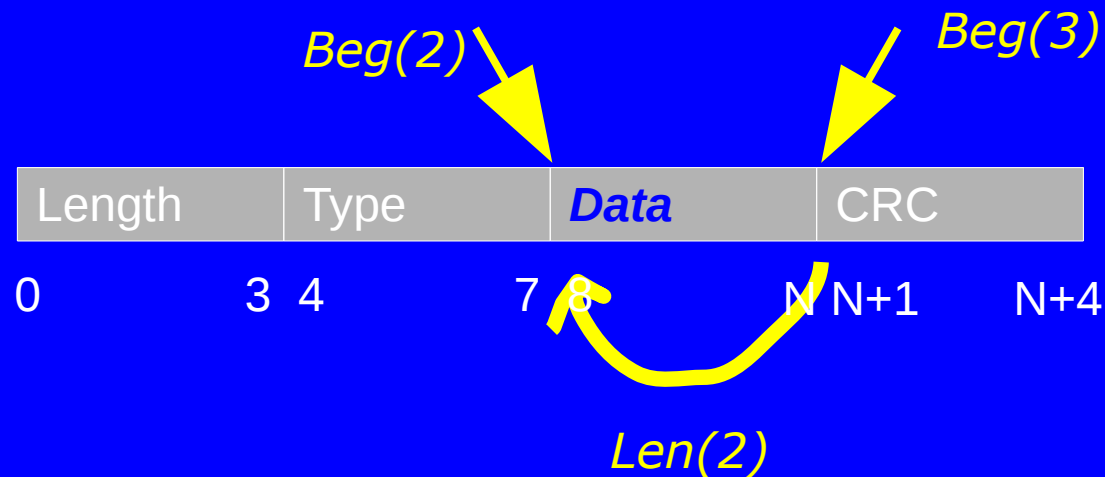


4d: Formalize parser's behaviour

- $Beg(i) \wedge Beg(i+1) \Rightarrow Len(i)$

join

- Compute the length of a field.

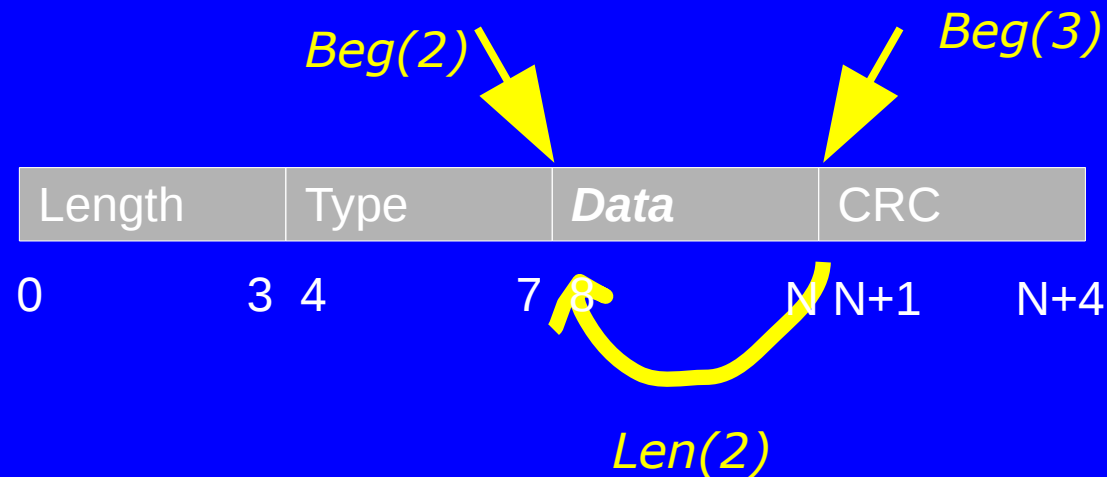


4d: Formalize parser's behaviour

- $Beg(i) \wedge Beg(i+1) \Rightarrow Len(i)$

join

- Compute the length of a field.



Deserializability Check Algorithm

- Transform a layout into a Horn KB: $O(3n)$
- Apply forward chaining: $O(3n)$
- Is $Len(i)$ for all i in a layout in KB? $O(n)$

Yes: Layout is deserializable

No: Layout is not deserializable.

Implementation

$(\text{Length} \rightarrow 4)\text{Length } f \ v \ f$

$(\text{Length} \rightarrow 4)\text{Length}_0 \ f_1 \ v_2 \ f_3$

Axioms

CLIPS

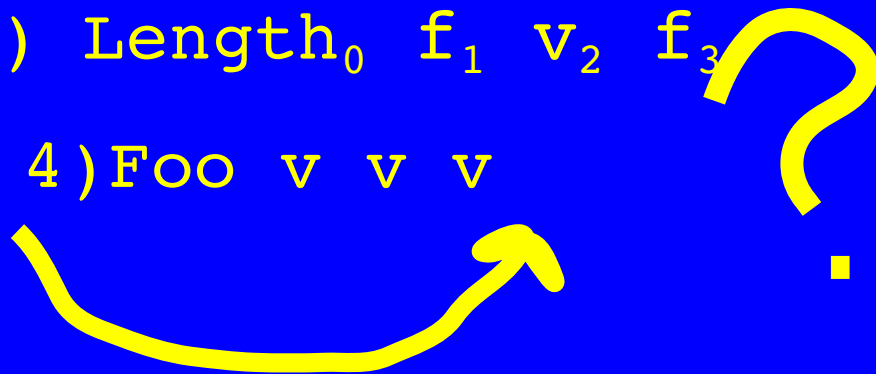
Python

$\forall i \Rightarrow$
 $\text{len}(i) ?$

Yes

Necessary condition for deserialization

- If layout L is deserializable, THEN in L
 - for every v_i
 - There is a $(foo \rightarrow s)_p, Xfoo_q$
 - Such that $q \leq i < q+s$
 - e.g. : $f \ v \ f$
 - $(Length \rightarrow 4) \ Length_0 \ f_1 \ v_2 \ f_3$
 - But: $(Foo \rightarrow 4) \ Foo \ v \ v \ v$



Repetition (Kleene star) : []*

(Foo → 2) Foo f [f v f]*

Repetition (Kleene star) : []*

(Foo → 2) Foo f [f v f]*

(Foo → 2) Foo₀ f₁ [f_{2.0} v_{2.1} f_{2.2}]*₂

**List labels
instead of natural labels**

[]*: Predicates

$(\text{Foo} \rightarrow 2) \text{Foo}_0 \text{f}_1 [\text{f}_{2.0} \text{v}_{2.1} \text{f}_{2.2}]^*_2$

- *Rep(i,l)* : field *i* is a repetition containing *l* fields
- *Replen(i)* : the parser knows repetition *i*'s length

What about the axioms?

- Lifted to the list label level:
- e.g:

$$Beg(b.a) \wedge Beg(b.a+1) \Rightarrow Len(b.a) \quad \text{join}$$

$$Ptr(b.a,s,i) \wedge Val(i) \wedge Beg(b.a) \quad \text{Jump right} \\ \Rightarrow Beg(b.a+s)$$

b,i : list

s,a : natural

[]*: Axioms

$(\text{Foo} \rightarrow 2) \text{Foo}_0 \text{f}_1 [\text{f}_{2.0} \vee_{2.1} \text{f}_{2.2}]^*_2$

- $\text{True} \Rightarrow \text{Rep}(i,l)$

in this layout: $\text{True} \Rightarrow \text{Rep}(2,3)$

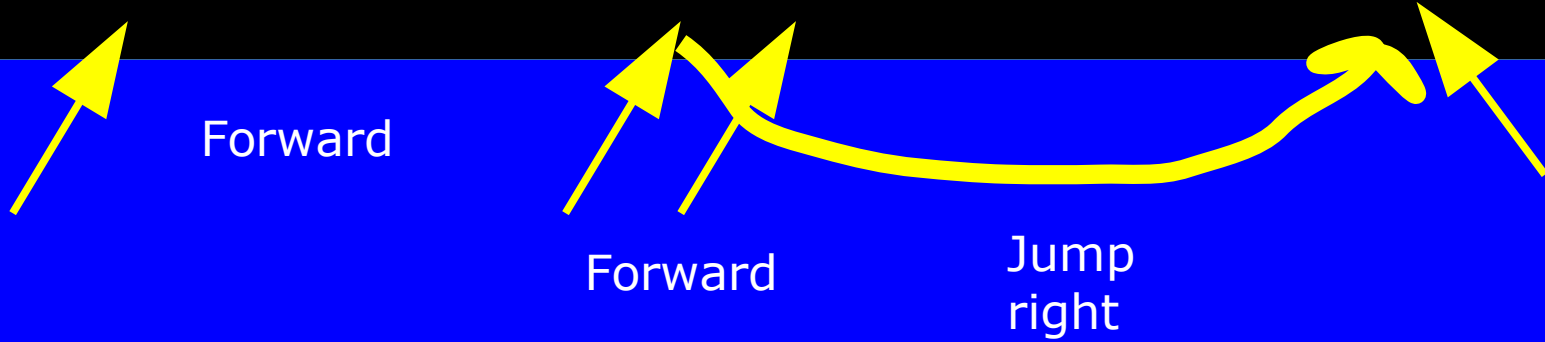
- $\text{Rep}(b.a,l) \wedge \text{Beg}(b.a) \wedge \text{Beg}(b.a+1) \Rightarrow \text{RepLen}(b.a)$

- $\text{Rep}(b.a,l) \wedge \text{Beg}(b.a) \Rightarrow \text{Beg}(b.a.0)$

- $\text{Rep}(b.a,l) \wedge \text{Beg}(b.a+1) \Rightarrow \text{Beg}(b.a.l)$

PITFALL!

$(\text{Foo} \rightarrow 2)\text{Foo}_0 \text{f}_1 [\text{f}_{2.0} \text{v}_{2.1} \text{f}_{2.2}]^*_2$



PITFALL!

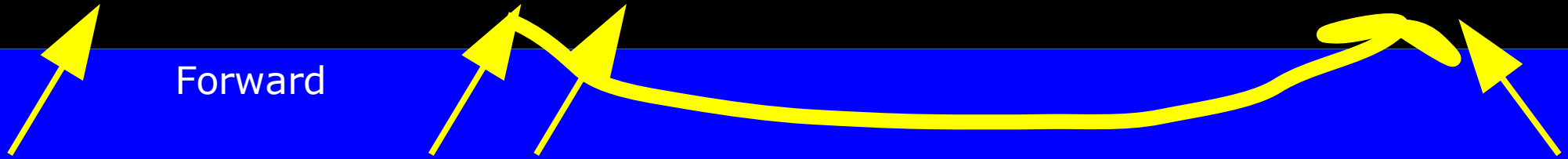
$(\text{Foo} \rightarrow 2)\text{Foo}_0 \ f_1 \ [f_{2.0} \ v_{2.1} \ f_{2.2} \]^*_2$



- $\text{Rep}(2,3) \wedge \text{Beg}(2) \Rightarrow \text{Beg}(2.0)$
- $\text{Rep}(2,3) \wedge \text{Beg}(3) \Rightarrow \text{Beg}(2.3)$

PITFALL!

`(Foo → 2)Foo0 f1 [f2.0 v2.1 f2.2 f2.3 v2.4 f2.5]*2`



Forward

Forward

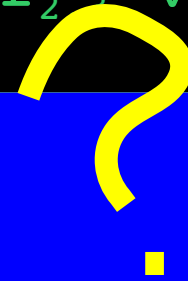
Jump
forward

`(Foo → 2)Foo0 f1 [f2.0 v2.1 f2.2]*2`

PITFALL!

$(\text{Foo} \rightarrow 2)\text{Foo}_0 \ f_1 \ [f_{2.0} \ v_{2.1} \ f_{2.2} \ f_{2.3} \ v_{2.4} \ f_{2.5} \]^*_2$

forward



backward

- $\text{Rep}(2,3) \wedge \text{Beg}(2) \Rightarrow \text{Beg}(2.0)$
- $\text{Rep}(2,3) \wedge \text{Beg}(3) \Rightarrow \text{Beg}(2.3)$

$(\text{Foo} \rightarrow 2)\text{Foo}_0 \ f_1 \ [f_{2.0} \ v_{2.1} \ f_{2.2} \]^*_2$

Dirty trick

$(\text{Foo} \rightarrow 2)\text{Foo}_0 \text{f}_1 [\text{f}_{2.0} \text{v}_{2.1} \text{f}_{2.2}]^*_2$

$(\text{Foo} \rightarrow 2)\text{Foo}_0 \text{f}_1 [\text{f}_{2.0} \text{v}_{2.1} \text{f}_{2.2} \text{f}_{2.3} \text{v}_{2.4} \text{f}_{2.5}]^*_2$

Take each repetition field and double its content.

If L2 is (not) deserializable
 Then all LN>2 are (not) deserializable

formally guaranteed:
 If L2 is deserializable,
 Then L is deserializable

L $(A \rightarrow 2)A_0 \ f_1 \ [f_{2.0}$

$[(B \rightarrow 1)B, \ v]^*_{2.1}$

$f_{2.2}]^*_2$

$(A \rightarrow 2)A_0 \ f_1 \ [f_{2.0}$

L2

$[(B \rightarrow 1)B \ v \ (C \rightarrow 1)C \ v]^*_{2.1}$

$f_{2.2} \ f_{2.3} \ [D \rightarrow 1)D \ v]^*_{2.4} \ f_{2.5}]^*_2$

$O(kn)$

Axioms are left undisturbed.

Implementation

$(\text{Length} \rightarrow 4)\text{Length} [f v]^* f$

$(\text{Length} \rightarrow 4)\text{Length} [f v f v]^* f$

$(\text{Length} \rightarrow 4)\text{Length}_0 [f_{2.0} v_{2.1} f_{2.2} v_{2.3}]^*_2 f_3$

Axioms



Python

$\forall i \Rightarrow \text{len}(i) ?$

$\forall \text{Rep}(i,j) \Rightarrow \text{RepLen}(i) ?$

NO

Intended application areas

- Serialization libraries
- Data definition language C!C
 - Rule-based parser generation?
 - Associate to each proof of deserializability a parser.

Related work

- Erlang, haskell, c...
- Pads
- Protocol buffers, avro, cap'n'proto, bson...

Summary

- Axiomatization of left-to-right stream parsing
- Implementation: Python+CLIPS
- Interesting results:
 - Necessary condition for deserializability
 - Doubling repetitions

Gabriele Paganelli

<https://github.com/gapag/horn-binary-deserialization>

<https://gapag.noblogs.org/>

gapag@distruzione.org