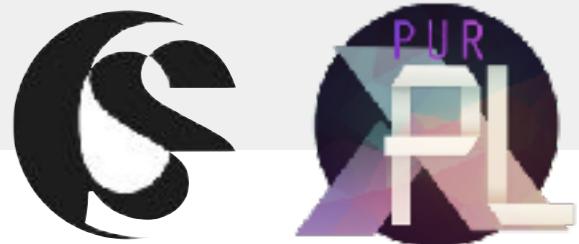


# NARCISSUS: Correct-By-Construction Derivation of Decoders and Encoders from Binary Formats

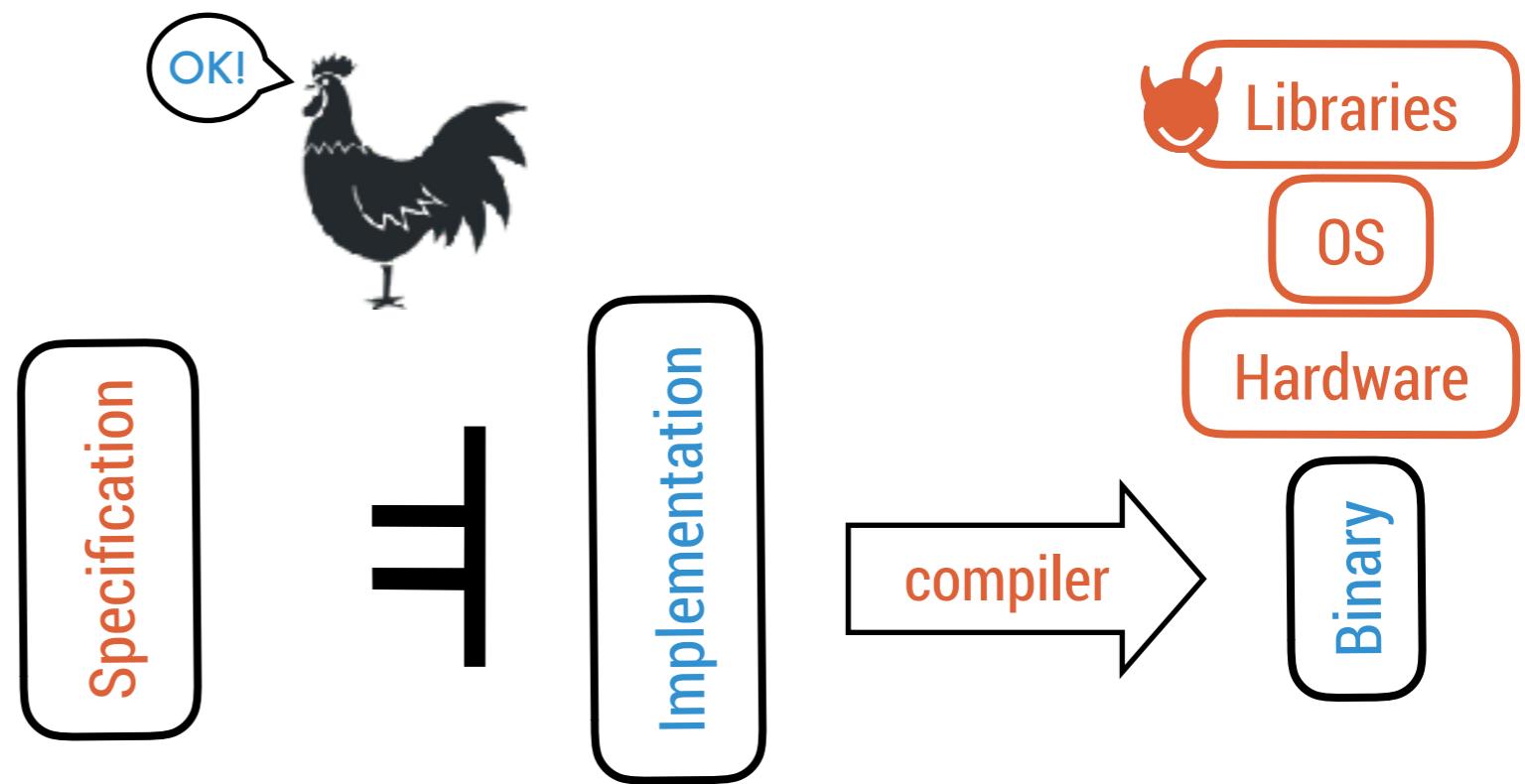
Benjamin Delaware  
Purdue University

Sorawit Suriyakarn<sup>1</sup>, Clément Pit-Claudel<sup>1</sup>,  
Qianchuan Ye<sup>2</sup>, and Adam Chlipala<sup>1</sup>  
<sup>1</sup>MIT, <sup>2</sup>Purdue University



# The Larger Picture

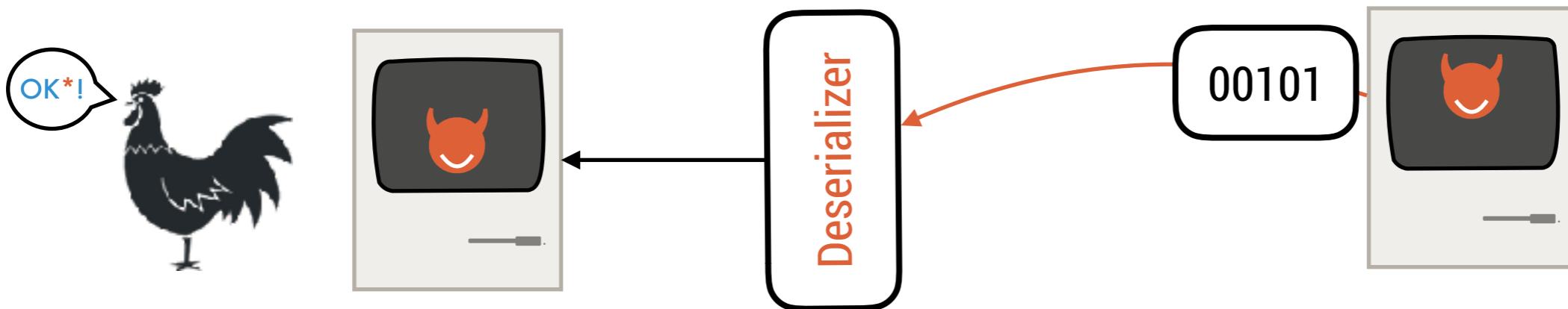
- Numerous developments of high-assurance software in proof assistants in the past five years:
  - CompCert C compiler
  - CertikOS
  - FSCQ file system
- Assurance comes from formal guarantees\* provided by proof assistant:



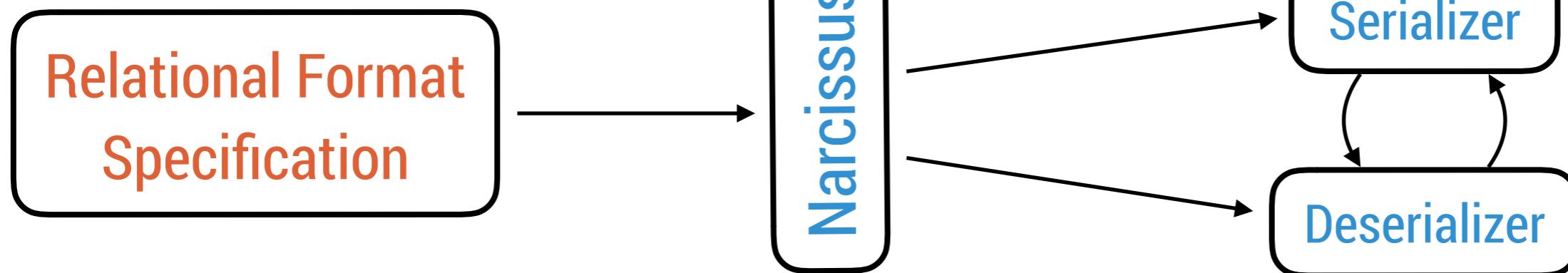
\* w.r.t Trusted Base

# Narcissus

- For networked systems, deserialization is important<sup>1</sup>
  - If these are in the TCB, bugs will break the assurance case!



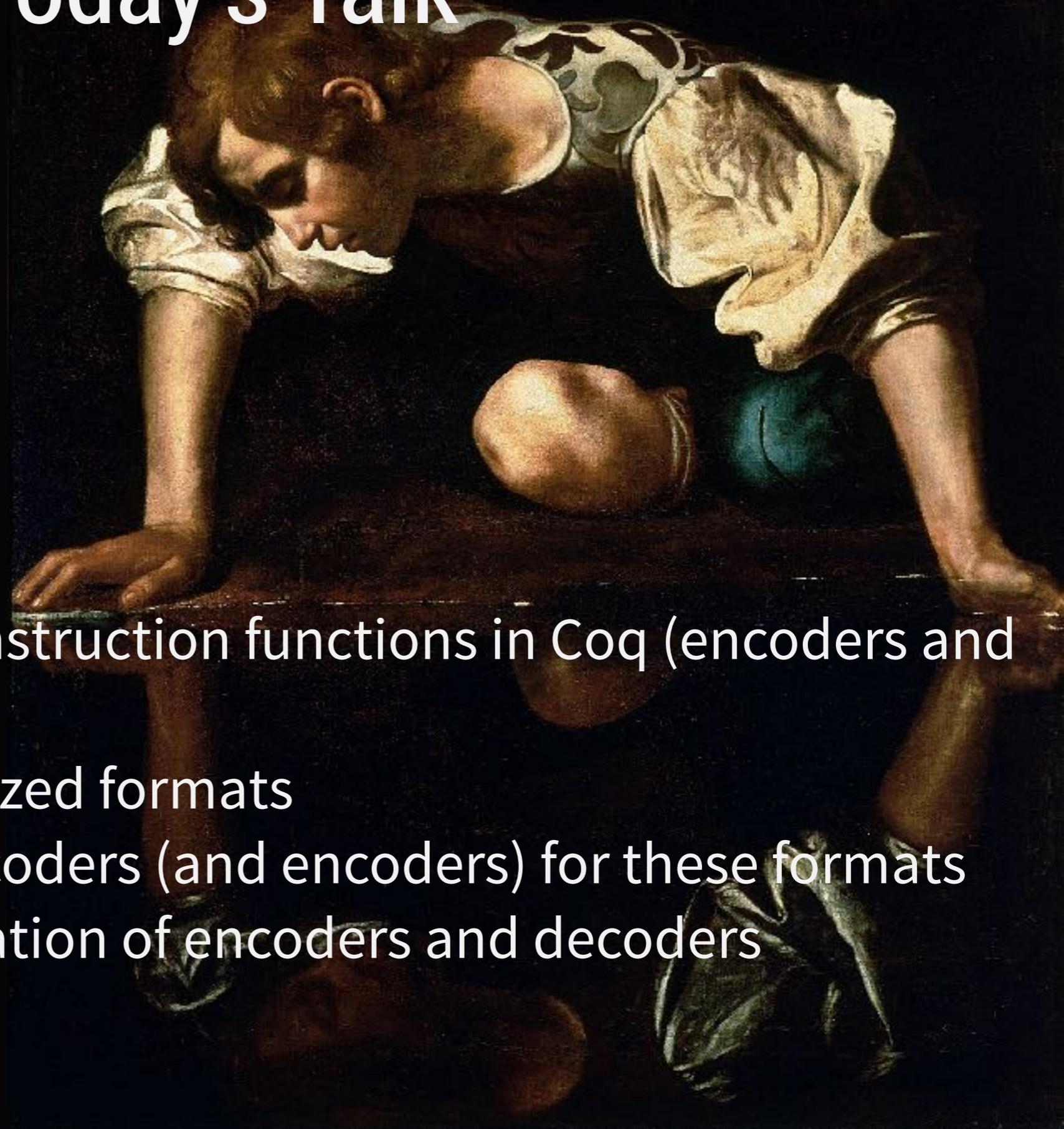
- Enter Narcissus:
  - User-extensible framework for deriving encoders and decoders from format specifications, with machine-checked correctness proofs



[1] An Empirical Study on the Correctness of Formally Verified Distributed Systems. Pedro Fonseca, Kaiyuan Zhang, Xi Wang, and Arvind Krishnamurthy. 2017.

# Today's Talk

- Deriving correct-by-construction functions in Coq (encoders and decoders)
- Formalizing standardized formats
- Specifying correct decoders (and encoders) for these formats
- Automating the derivation of encoders and decoders



# A Simple Format

## SOURCE TYPE

**Definition** SimpleProd n :=  
Word (\* id \*)  
× DateTime (\* Timestamp \*)  
× Vector Word n. (\* Payload \*)

## FORMAT

**Definition** SimpleProd\_fmt n  
: FormatM (SimpleProd n) T :=  
Pair\_fmt  
Word\_fmt  
(Pair\_fmt  
DTime\_fmt  
(Vector\_fmt Word\_fmt)).

## RFC-STYLE FORMAT

0	1	2	3	4	5	6	7
		ID					
		TIMESTAMP					
/		PAYLOAD					/
/							/

## DECODER

**Definition** SimpleProd\_dec n  
: DecodeM (SimpleProd n) T :=  
Pair\_dec  
Word\_dec  
(Pair\_dec  
DTime\_dec  
(Vector\_dec Word\_dec)).

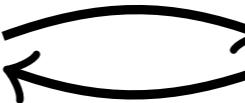
- SimpleProd\_fmt = composition of sub-formats

• Prod\_fmt :  $\forall S_1 S_2, \text{Format } S_1 T \rightarrow \text{Format } S_2 T \rightarrow \text{FormatM } (S_1 \times S_2) T$

- SimpleProd\_dec = composition of decoders

# Correct-by-Construction Decoders

- Format and decoder is so similar, let's define a combinator<sup>2</sup> packaging the two together **with a proof of correctness**:

```
Record CorrectFormatDecoder' (S T : Type) :=  
{ format : FormatM S T; decode : DecodeM S T;  
decodeOK : format  decode }
```

“forall inputs, decode **correctly inverts** format”

- Correct decoder = inhabitant of this type
- Coq is quite good at finding inhabitants of types!
- **Deductive Synthesis:** “To construct a program whose output satisfies the conditions of the specification, we prove a theorem stating the existence of such an output. The proof is restricted to be sufficiently constructive so that a program computing the desired output can be extracted directly from the proof. The proof constitutes a demonstration of the correctness of this program<sup>3</sup>.”

[2] FUNCTIONAL PEARL Pickler Combinators. Andrew J. Kennedy. 2004.

[3] Deductive Synthesis Of The Unification Algorithm. Richard J. Waldinger. 1981.

# Decoder Combinators

- Proof-Carrying combinators are also compositional:

```
Definition ProdComb (S1 S2 T : Type)
  (fmt1 : CorrectFormatDecoder' S1 T) (fmt2 : CorrectFormatDecoder' S2 T)
  : CorrectFormatDecoder' (S1 * S2) T :=
{ format := Pair_fmt fmt1.format fmt2.format;
  decode := Pair_dec fmt1.decode fmt2.decode;
  decodeOK := Pair_correct fmt1.decodeOK fmt2.decodeOK }
```

- With some combinators, we can define interactively:

```
Definition SimpleRec_Correct n :
  CorrectFormatDecoder' (SimpleRec n) T.
```

**Proof.**

```
apply ProdComb.
- apply WordComb.
- apply ProdComb.
  + apply DateTimeComb.
  + apply (VectorComb n).
```

**Defined.**

## FORMAT

```
Definition SimpleProd_fmt n
  : FormatM (SimpleProd n) T :=
Pair_fmt
Word_fmt
(Pair_fmt
DTime_fmt
(Vector_fmt Word_fmt)).
```

# Decoder Combinators

- Proof-Carrying combinators are also compositional:

```
Definition ProdComb (S1 S2 T : Type)
  (fmt1 : CorrectFormatDecoder' S1 T) (fmt2 : CorrectFormatDecoder' S2 T)
  : CorrectFormatDecoder' (S1 * S2) T :=
{ format := Pair_fmt fmt1.format fmt2.format;
  decode := Pair_dec fmt1.decode fmt2.decode;
  decodeOK := Pair_correct fmt1.decodeOK fmt2.decodeOK }
```

- With some combinators, we can define interactively:

```
Definition SimpleRec_Correct n :
  CorrectFormatDecoder' (SimpleRec n) T.
```

**Proof.**

**auto**.

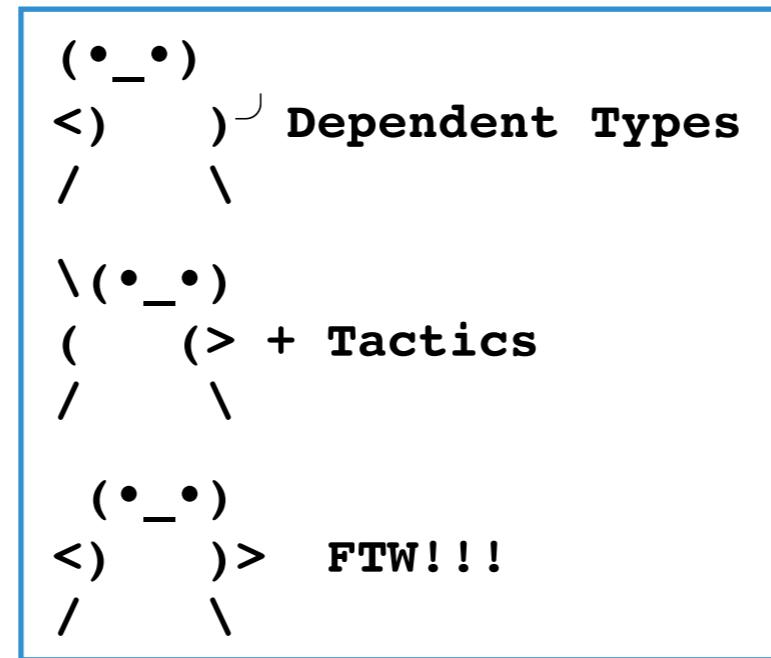
**Defined.**

## FORMAT

```
Definition SimpleProd_fmt n
  : FormatM (SimpleProd n) T :=
  Pair_fmt
    Word_fmt
    (Pair_fmt
      DTime_fmt
      (Vector_fmt Word_fmt)).
```

# Recap

- We're done\*:
  - Define a package of formats + decoders together with proofs
  - Define compositional instances of these records
  - Use tactics to automate derive implementations automatically.



- Well, **maybe not**. This format made several simplifying assumptions that don't apply to many common formats.

\*modulo some elided definitions

# A Not as Simple Format

- What about this example?

- Source Type:

- Record Type
- Variable-length payload

- Format:

- Encodes payload length
- Out of sync with type:
  - $\text{++} : \forall S T, \text{FormatM } S T \rightarrow \text{FormatM } S T \rightarrow \text{FormatM } S T$
  - $\circ : \forall S S' T, \text{FormatM } S' T \rightarrow (S \rightarrow S') \rightarrow \text{FormatM } S T$

- Decoder:

- Uses decoded length to parse payload
- Needs to reconstruct source value

SOURCE TYPE

```
Record SimpleRec :=  
  { id : Word;  
    Payload : List Word;  
    Timestamp : DateTime }.
```

FORMAT

```
Definition SimpleRec_fmt  
  : FormatM SimpleRec T :=  
    Word_fmt o id  
  + Nat_fmt 8 o length o Payload  
  + DTime_fmt o Timestamp  
  + List_fmt Word_fmt o Payload.
```

DECODER

```
Definition SimpleRec_dec  
  : DecodeM SimpleRec T :=  
    w ← Word_dec;  
    ln ← Nat_dec;  
    dt ← DTime_dec;  
    l ← List_dec Word_dec ln;  
    return { id := w; Payload := l;  
           Timestamp := ln }.
```

# Decoding Lists

- How to define proof-carrying combinator for Lists?

```
Definition List_Correct' {S T} (len : nat) fmts
  (fmtsOK : CorrectFormatDecoder' S T fmts)
    : CorrectFormatDecoder' (List S) T (List_fmt fmts) :=
{ decode := List_dec fmts.decode len;
  decodeOK := ?? len }
```

- $\text{??} : \text{List\_fmt } \text{fmts} \xleftrightarrow{\hspace{1cm}} \text{List\_dec } \text{fmts.decode } n$
- $\text{List\_dec } \text{fmts.decode } len$  builds lists of length  $len$ , but the intuitive spec talks about all source values!
- Correctness of the decoder depends on the larger format
- Goal: Compositional, context-aware decoder combinators

# Constraining the Context

- Proposal: Restrict the source values considered:

format|<sub>P</sub>  decode

- “For all source values  $s$  such that  $P s$ , decode correctly inverts format”
- Can now define a combinator for lists:

```
Definition List_Correct' {S T} (len : nat) fmts
  (fmtsOK : CorrectFormatDecoder' S T fmts)
    : CorrectFormatDecoder' (List S) T (List_fmt fmts) | {l | length l = len}
  := {decode := List_dec fmts.decode len;
     decodeOK := list_correct len }
```

- Works for other dependencies: e.g. tags, version numbers, etc.

# Threading Context

- How to ensure that proof-relevant context gets passed down?

```
Definition Seq_Word_Correct {S T} P fmtS (f : S -> Word)
  (fmtSOK : ∀ w, CorrectFormatDecoder'' S T fmtS | (P n { s | f s = w })
  : CorrectFormatDecoder'' S T (Word_fmt ∘ f ++ fmtS) |P :=
{ decode := Word_dec >>= fmtS.decode;
  decodeOK := seq_word_OK fmtS f fmtSOK.decodeOK
```

- Finding a correct decoder for a sequence of formats only needs to consider the case where the first format was decoded correctly

$P \cap \{ s \mid f s = w \}$

“previous context”

“newly decoded context”

# Threading Context

CorrectFormatDecoder  $\left( \begin{array}{l} \text{Word\_fmt } \circ \text{id} \\ \# \text{Nat\_fmt } 8 \circ \text{length } \circ \text{Payload} \\ \# \text{DTime\_fmt } \circ \text{Timestamp} \\ \# \text{List\_fmt Word\_fmt } \circ \text{Payload} \end{array} \right) \quad \{s \mid \text{length s.Payload} < 2^8\}$

CorrectFormatDecoder  $\left( \begin{array}{l} \text{Nat\_fmt } 8 \circ \text{length } \circ \text{Payload} \\ \# \text{DTime\_fmt } \circ \text{Timestamp} \\ \# \text{List\_fmt Word\_fmt } \circ \text{Payload} \end{array} \right) \quad \{s \mid \text{length s.Payload} < 2^8\} \cap \{s \mid s.id = w\}$

CorrectFormatDecoder  $\left( \begin{array}{l} \text{DTime\_fmt } \circ \text{Timestamp} \\ \# \text{List\_fmt Word\_fmt } \circ \text{Payload} \end{array} \right) \quad \{s \mid \text{length s.Payload} < 2^8\} \cap \{s \mid s.id = w\} \cap \{s \mid \mid s.Payload \mid = ln\}$

CorrectFormatDecoder  $\left( \begin{array}{l} \text{List\_fmt Word\_fmt } \circ \text{Payload} \end{array} \right) \quad \{s \mid \text{length s.Payload} < 2^8\} \cap \{s \mid s.id = w\} \cap \{s \mid \mid s.Payload \mid = ln\} \cap \{s \mid s.Timestamp = dt\}$

CorrectFormatDecoder  $\epsilon \quad \{s \mid \text{length s.Payload} < 2^8\} \cap \{s \mid s.id = w\} \cap \{s \mid \mid s.Payload \mid = ln\} \cap \{s \mid s.Timestamp = dt\} \cap \{s \mid s.Payload = l\}$

# Threading Context

- The last step is constrained to a single value
- It's obvious how to rebuild the original source value

CorrectFormatDecoder     $\epsilon \uparrow$

$$\begin{aligned} & \{s \mid \text{length } s.\text{Payload} < 2^8\} \\ & \cap \{s \mid s.\text{id} = w\} \\ & \cap \{s \mid |s.\text{Payload}| = ln\} \\ & \cap \{s \mid s.\text{Timestamp} = dt\} \\ & \cap \{s \mid s.\text{Payload} = l\} \end{aligned}$$

## DECODER

**Definition** SimpleRec\_dec  
: DecodeM SimpleRec T :=  
w  $\leftarrow$  Word\_dec;  
ln  $\leftarrow$  Nat\_dec;  
dt  $\leftarrow$  DTime\_dec;  
l  $\leftarrow$  List\_dec Word\_dec ln;  
**return** { id := w; Payload := l;  
Timestamp := dt }.

# Base Combinators

- Our basic building blocks:
  - Functional programs in monadic style
  - Target type is a monoid
    - `snoc t x = t ++ [x]`
    - `unfold [x] ++ t = (x, t)`
- Narcissus-specific details hidden in bind operation
- Proofs components are proofs

## FORMAT

```
Fixpoint Word_fmt {n}
  : FormatM (Word n) T :=
  λw. match w with
    | ε ⇒ t
    | b · w' ⇒ t ← Word_fmt w';
      return (snoc b t)
  end.
```

## DECODER

```
Fixpoint Word_dec {n}
  : DecodeM (Word n) T :=
  λt. match n with
    | O ⇒ return ε
    | S n' ⇒ (b, t') ← unfold t;
      w' ← word_dec n' t';
      return (w' · b)
  end.
```

# Nondeterministic Formats

- Many formats do not prescribe a canonical target
  - Version = 8-bit word whose two lowest-order bits are set
  - $2^6$  different encodings for each source value
- Many Examples: Ethernet, ASN.1, DNS
- Narcissus' solution: relational interpretation of formats
  - Every source value is related to a set of target values
  - Admits new choice operator:  $\{x \mid P x\}$

RFC-STYLE FORMAT							
0	1	2	3	4	5	6	7
+	+	+	+	+	+	+	+
			ID				
+	+	+	+	+	+	+	+
			LENGTH				
+	+	+	+	+	+	+	+
	1	1		VERSION			
+	+	+	+	+	+	+	+
			TIMESTAMP				
+	+	+	+	+	+	+	+
/			PAYLOAD				/
/							/
+	+	+	+	+	+	+	+

FORMAT

**Definition** NDRec\_fmt

```
: FormatM SimpleRec T :=  
  Word_fmt o id  
  # Word_fmt o length o Payload  
    (* New: Version Number *)  
  # Word_fmt o {w | w ≥ 11000000}  
  # DTime_fmt o Timestamp  
  # Vector_fmt Word_fmt o Payload.
```

# IPv4

- Extensions = User-Defined Components!

```

Definition IPv4_Packet_Format : FormatM IPv4_Packet T :=
  format_word ○ (fun _ => natToWord 4 4)
++ format_nat 4 ○ IPv4_Packet_Header_Len
++ format_unused_word 8 (* TOS Field! *)
++ format_word ○ TotalLength
++ format_word ○ ID
++ format_unused_word 1 (* Unused flag! *)
++ format_bool ○ DF
++ format_bool ○ MF
++ format_word ○ FragmentOffset
++ format_word ○ TTL
++ format_enum ProtocolTypeCodes ○ Protocol
++ IPChecksum ++
++ format_word ○ SourceAddress
++ format_word ○ DestAddress
++ format_list format_word ○ Options.

```

# The Emperor's Old Clothes

Charles Antony Richard Hoare  
Oxford University, England

I gave desperate warnings against the obscurity, the complexity, and overambition of the new design, but my warnings went unheeded. I conclude that there are two ways of constructing a software design: One way is to make it so simple that there are obviously no deficiencies and the other way is to make it so complicated that there are no obvious deficiencies.

## A Constructive Approach to the Problem of Program Correctness.

Personally I cannot refrain from feeling that many debugging aids that are en vogue now are invented as a compensation for the shortcomings of a programming technique that will be denounced as obsolete within the near future.

# Design Decisions

- Nondeterministic, relational formats, with Narcissus-specific details hidden in monads
- Combinator libraries for decoders (and encoders) keyed on a common format.
- Automate derivations when possible using tactics
- Dependencies between subformats are reflected in the proofs of correctness for individual combinators.

On to the details!

# Specifying Formats

- Type of Formats:
  - $\text{FormatM } S \ T \ \Sigma := \text{Set of}(S \times \Sigma \times T \times \Sigma)$
- Define relations via Nondeterminism monad:
  - $e \in \text{return } v \leftrightarrow e = v$
  - $e \in \{ x \mid P x \} \leftrightarrow P e$
  - $e \in x \leftarrow y; \ k \ x \leftrightarrow \exists e'. e' \in y \wedge e \in k \ e'$
- Or using higher-order formats

## SEQUENCE FORMATS

$$(s, \sigma, t, \sigma') \in \text{format}_1 + \text{format}_2 \leftrightarrow \exists t_1 t_2 \sigma''. (s, \sigma, t_1, \sigma'') \in \text{format}_1 \wedge (s, \sigma'', t_2, \sigma') \in \text{format}_2 \wedge t = t_1 \cdot t_2$$

## COMPOSE FORMATS

$$(s, \sigma, t, \sigma') \in \text{format} \odot f \leftrightarrow \exists s'. (s', \sigma, t, \sigma') \in \text{format} \wedge f s s'$$

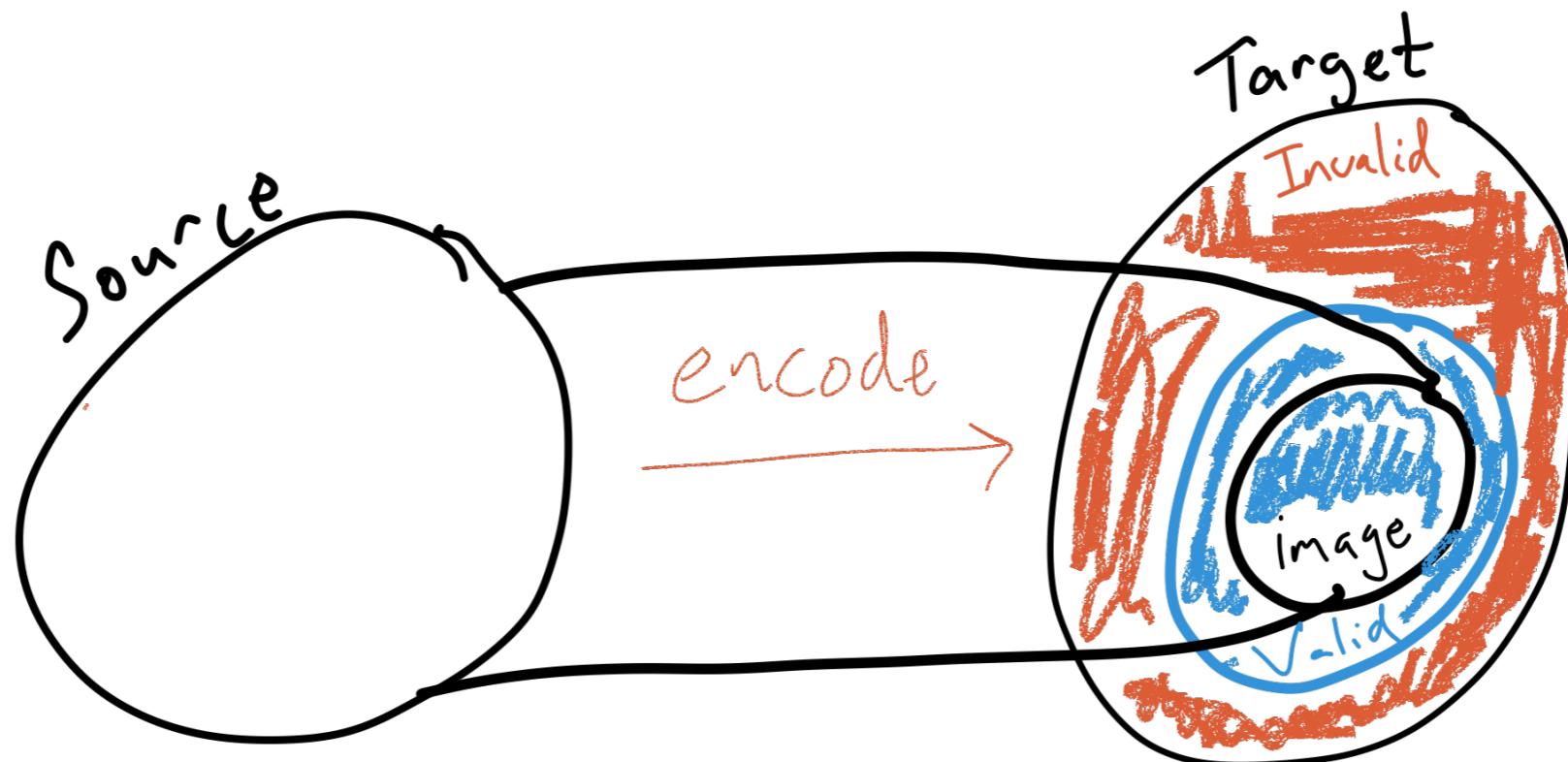
## UNION FORMAT

$$(s, \sigma, t, \sigma') \in \bigcup_i \text{format}_i \leftrightarrow \exists j. (s, \sigma, t, \sigma') \in \text{format}_j$$

# Correct Encoders

- A correct **encoder** is a function wholly contained in the relation defined by a format

$\text{EncodeM } S \ T \ \Sigma := S \rightarrow \Sigma \rightarrow \text{Option } T \times \Sigma$



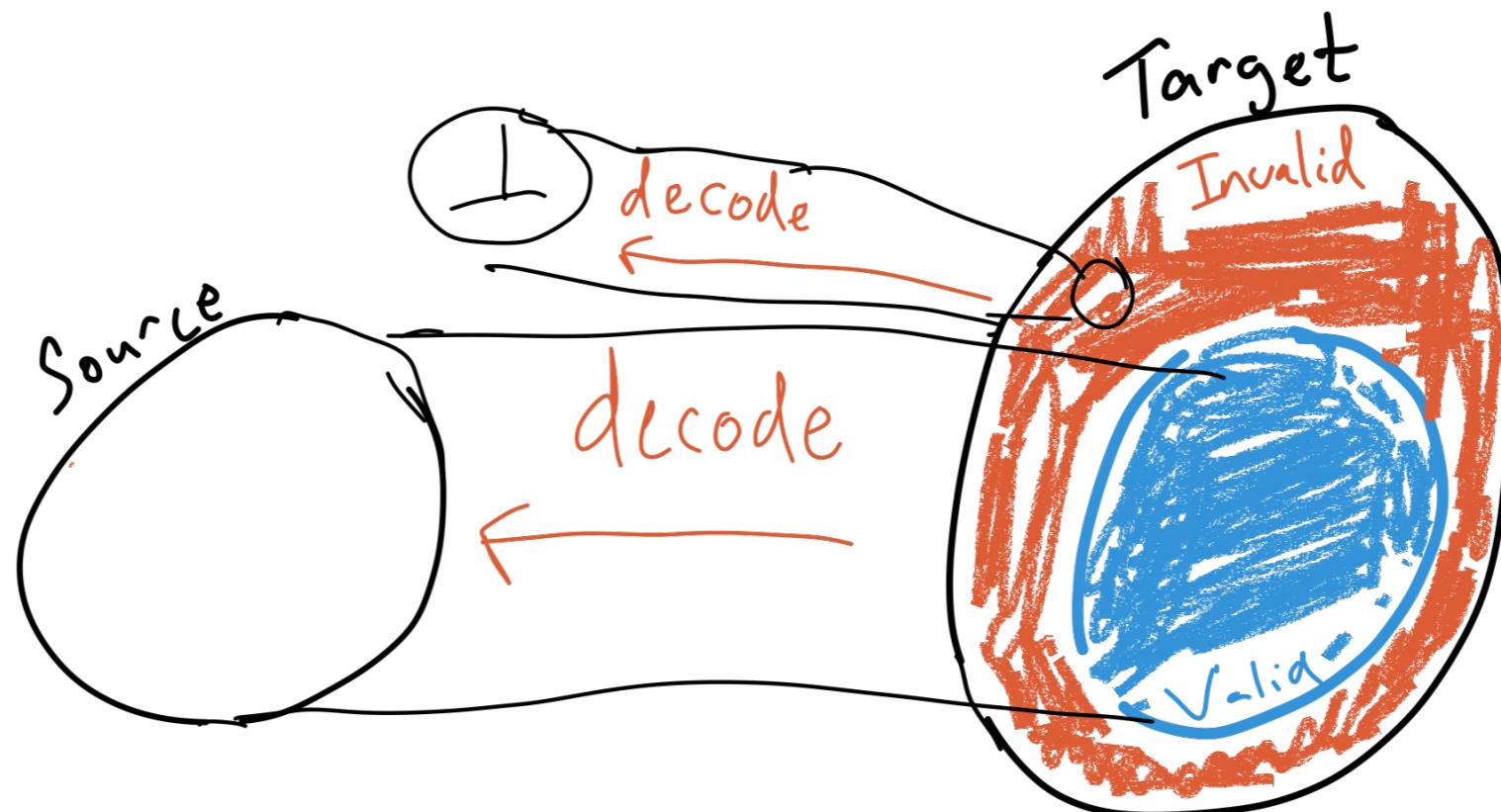
$\text{encode} \subseteq \text{format} :=$

$\forall s \ \sigma \ t \ \sigma'. \text{encode } s \ \sigma = \text{Some } (t, \sigma') \rightarrow (s, \sigma, t, \sigma') \in \text{format}$

$\wedge \forall s. \text{encode } s \ \sigma = \text{None} \rightarrow \forall t' \ \sigma'. (s, \sigma, t', \sigma') \notin \text{format}$

# Correct Decoders

- A correct **decoder** maps values in the image of the format back to the original source value, *and* signals an error for other values



format  $\xleftrightarrow{\approx}$  decode.

$$\left( \begin{array}{l} \forall (\sigma_E : \Sigma_E) (\sigma_D : \Sigma_D) (s : S) (t : T). \\ (s, \sigma_E, t, \sigma_D) \in \text{format} \wedge \sigma_E \approx \sigma_D \\ \rightarrow \exists \sigma_D'. \text{decode } t \sigma_D = \text{Some } (s, \sigma_D') \\ \wedge \sigma_E' \approx \sigma_D' \end{array} \right) \quad \text{ii} \quad \left( \begin{array}{l} \forall (\sigma_E : \Sigma_E) (\sigma_D : \Sigma_D) (s : S) (t : T). \\ \sigma_E \approx \sigma_D \wedge \text{decode } t \sigma_D = \text{Some } (s, \sigma_D') \\ \rightarrow \exists \sigma_E'. (s, \sigma_E, t, \sigma_E') \in \text{format} \\ \wedge \sigma_E' \approx \sigma_D' \end{array} \right)$$

# Relating Encoders + Decoders

- **Encoders** and **decoder** are specified by a common **format**:

encode  $\subseteq$  format  $\xleftarrow[\approx]{}$  decode.

- Correctness of both yields the expected roundtrip properties:

**THEOREM 2.3 (DECODE INVERTS ENCODE).** *Given a correct decoder format  $\xleftarrow[\approx]{}$  decode and correct encoder encode  $\subseteq$  format for a common format format, decode is an inverse for encode when restricted to source values in the format:*

$$\forall s \sigma_E t \sigma'_E \sigma_D. \text{encode } s \sigma_E = \text{Some}(t, \sigma'_E) \wedge \sigma_E \approx \sigma_D \rightarrow \exists \sigma'_D. \text{decode } t \sigma_D = \text{Some}(s, \sigma'_D)$$

**THEOREM 2.4 (ENCODE INVERTS DECODE).** *Given a correct decoder format  $\xleftarrow[\approx]{}$  decode and correct encoder encode  $\subseteq$  format for a common format format, encode is defined for all decoded source values produced by decode,*

$$\forall s \sigma_D t \sigma'_D \sigma_E. \text{decode } t \sigma_D = \text{Some}(s, \sigma'_D) \wedge \sigma_E \approx \sigma_D \rightarrow \exists t' \sigma'_E. \text{encode } s \sigma_E = \text{Some}(t', \sigma'_E)$$

# Encoder Combinators

- Combinators for correct-by-construction encoders are similar to decoders:

```
Record CorrectFormatEncoder (S T Σ : Type) :=
  { format : FormatM S T Σ ; encode : EncodeM S T Σ ;
    encodeOK : encode ⊆ format }
```

- Formats define a family of valid encoding strategies;
- A correct encoder picks one:

```
Fixpoint Word_enc {n} {σ}
  : EncodeM (Word n) T σ:=
  λw. match w with
    | ε => t
    | b · w' => t ← Word_enc w';
      return (snoc b t)
  end.
```

```
Fixpoint Word_fmt {n} {σ}
  : FormatM (Word n) T σ:=
  λw. match w with
    | ε => t
    | b · w' => t ← Word_fmt w';
      return (snoc b t)
  end. (WORDENC)
```

# Deriving Encoders

- Encoders derivations pick a particular family member

$$\lambda s. t_1 \leftarrow \text{Word\_enc } s.\text{id};$$
$$t_2 \leftarrow \square;$$
$$\mathbf{return} (t_1 \cdot t_2) \subseteq \begin{array}{l} \text{Word\_fmt } \circ \text{id} \\ \# \text{Nat\_fmt } 8 \circ \text{length } \circ \text{Payload} \\ \# \text{DTime\_fmt } \circ \text{Timestamp} \\ \# \text{List\_fmt Word\_fmt } \circ \text{Payload} \end{array}$$

↑WORDSEQENC

$$\lambda s. t_2 \leftarrow \text{Nat\_enc}$$
$$(\text{length } s.\text{Payload});$$
$$t_3 \leftarrow \square;$$
$$\mathbf{return} (t_2 \cdot t_3) \subseteq \begin{array}{l} \text{Nat\_fmt } 8 \circ \text{length } \circ \text{Payload} \\ \# \text{DTime\_fmt } \circ \text{Timestamp} \\ \# \text{List\_fmt Word\_fmt } \circ \text{Payload} \end{array}$$

↑NATSEQENC

$$\lambda s. t_3 \leftarrow \text{DTime\_enc}$$
$$s.\text{Timestamp};$$
$$t_4 \leftarrow \square;$$
$$\mathbf{return} (t_3 \cdot t_4) \subseteq \begin{array}{l} \text{DTime\_fmt } \circ \text{Timestamp} \\ \# \text{List\_fmt Word\_fmt } \circ \text{Payload} \end{array}$$

↑DTSEQENC

$$\lambda s. \mathbf{return} (\text{List\_enc }$$
$$\text{Word\_enc} \subseteq \text{List\_fmt Word\_fmt } \circ \text{Payload}$$
$$s.\text{Payload})$$
$$\lambda s. t_1 \leftarrow \text{Word\_enc } s.\text{id};$$
$$t_2 \leftarrow \text{Nat\_enc } (\text{length } s.\text{Payload});$$
$$t_3 \leftarrow \text{DTime\_enc } s.\text{Timestamp};$$
$$t_4 \leftarrow \text{List\_enc Word\_enc } s.\text{Payload};$$
$$\mathbf{return} (t_1 \cdot t_2 \cdot t_3 \cdot t_4)$$

# Decoder Combinators

- Decoders for subformats have a particular shape:
  - They're partial: return rest of target value for further parsing
  - Restricted source values
- Bake this into the type of combinators:

```
Record CorrectFormatDecoder {S T ΣE ΣD}  
(sourcePred : S -> Prop) ≈ (fmts : FormatM S T ΣE) :=  
{ decode : DecodeM (S * T) T ΣD;  
  decodeOK : fmts |sourcePred  decode }
```

# Decoder Combinators

- Can define combinators for our higher-order formats:

$$\frac{\begin{array}{c} P' \cap \text{format}_1 \xrightarrow{\approx} \text{decode}_1 \\ \forall s. P s \rightarrow P'(f s) \\ \hline \forall s'. P' s' \rightarrow \{s \mid f s = s' \wedge P s\} \cap \text{format}_2 \xrightarrow{\approx} \text{decode}_2 s' \end{array}}{P \cap \text{format}_1 \circ f + \text{format}_2 \xrightarrow{\approx} \text{decode}_1 + D \text{ decode}_2} \quad (\text{SEQDEC})$$

- The assumptions characterize how to **correctly** decode a sequence of two formats:
  - Decoder for first format
  - Restrictions on first format are satisfied
  - Decoder for second format, assuming the argument is a projection of original source value

# Decoder Combinators

- Can define combinators for our higher-order formats:

$$\frac{\begin{array}{c} \hline P \cap \text{format}_i \xrightarrow{\approx} \text{decode}_i \\ \hline \forall (s, \sigma, t, \sigma') \in \bigcup_i \text{format}_i. (s, \sigma, t, \sigma') \in \text{format}_{n(t)} \end{array}}{\bigcup_i \text{format}_i \xrightarrow{\approx} \lambda t. j \leftarrow n t; \text{decode}_j t}$$

(UNIONDEC)

- The assumptions characterize how to **correctly** decode the union of several formats:
  - Decoder for each variant
  - Index function  $n$ , which correctly determines format used

# Decoder Combinators

- Can define combinators for our higher-order formats:

$$\frac{\begin{array}{c} P' \cap \text{format} \xrightarrow{\approx} \text{decode} \\ \forall s. P s \rightarrow P'(f s) \quad \forall s. P s \rightarrow (g \circ f)s = s \\ \forall s. P'(f s) \rightarrow ((b \circ g \circ f)s = \text{true} \leftrightarrow P s) \end{array}}{\lambda t. (s, t') \leftarrow \text{decode } t; \begin{array}{l} \text{if } b \circ g s \text{ then} \\ \quad \text{return } (g s) \\ \text{else fail} \end{array} \text{ (TRMDEC)}}$$

- The assumptions characterize how to decode a “singleton” format:
  - Decoder for first format
  - Restrictions on first format are satisfied
  - Restriction uniquely determines source value
  - Decision procedure for source restriction

# Decoder Combinators

- Consider the final proof state for our simple record format:

$$\frac{\begin{array}{c} P' \cap \text{format} \xrightarrow{\approx} \text{decode} \\ \forall s. P s \rightarrow P'(f s) \quad \forall s. P s \rightarrow (gof)s = s \\ \forall s. P' (f s) \rightarrow ((bogof) s = \text{true} \leftrightarrow P s) \end{array}}{\lambda t. (s, t') \leftarrow \text{decode } t;}$$

$P \cap \text{format} \circ f \xrightarrow{\approx}$       **if**  $b \circ g s$  **then**  
   **return**  $(g s)$   
   **else** fail  
   (TrMDEC)

$$\frac{\forall w \ln dt. \ln < 2^8 \rightarrow \{s \mid s.\text{Timestamp} = dt \wedge \text{length } s.\text{Payload} = \ln \wedge s.\text{id} = w \wedge \text{length } s.\text{Payload} < 2^8\} \cap \text{List\_fmt Word\_fmt} \circ \text{Payload}}{\lambda t. (l, t') \leftarrow \text{List\_dec Word\_dec} \ln t; \quad \text{return } \{\text{id} := w; \text{Timestamp} := dt; \text{Payload} := l\}}$$

# Automating Derivations

- Tricky bit is finding appropriate  $g$  and  $b$
- Tactics to the rescue!
- $\square_g$  is an evar for a function which reconstructs the original source value.
- Why not use reflexivity at start?

$$\begin{array}{l} H_1 : \text{ln} < 2^8 \\ H_2 : s.\text{Timestamp} = dt \wedge \text{length } s.\text{Payload} = \text{ln} \\ \quad \wedge s.\text{id} = w \wedge \text{length } s.\text{Payload} < 2^8 \\ \hline \square_g s.\text{Payload} = s \end{array}$$

*Destruct the source value,  $s$ .*

$$\begin{array}{l} H_1 : \text{ln} < 2^8 \\ H_2 : x_2 = dt \wedge \text{length}(x_3) = \text{ln} \wedge x_1 = w \\ \quad \wedge \text{length}(x_3) < 2^8 \\ \hline \square_g x_3 = \{ \text{id} := x_1; \\ \quad \text{Timestamp} := x_2; \\ \quad \text{Payload} := x_3 \} \end{array}$$

*Substitute with equalities in source restriction.*

$$\begin{array}{l} H_1 : \text{ln} < 2^8 \\ H_2 : \text{ln} < 2^8 \\ \hline \square_g x_3 = \{ \text{id} := w; \\ \quad \text{Timestamp} := dt; \\ \quad \text{Payload} := x_3 \} \end{array}$$

*Variant of reflexivity solves the goal.  $\square$*

# Simplifying Specifications

- Narcissus includes a library of common formats
  - Base formats for single data types
  - Higher-order formats for compositional format definitions

Format	LoC	LoP	Higher-order
Sequencing (+)	7	164	Y
Termination (e)	1	28	N
Conditionals	25	204	Y
Booleans	4	24	N
Fixed-length Words	65	130	N
Unspecified Field	30	60	N
List with Encoded Length	40	90	N
String with Encoded Length	31	47	N
Option Type	5	79	N
Ascii Character	10	53	N
Enumerated Types	35	82	N
Variant Types	43	87	N
Domain Names	86	671	N
IP Checksums	15	1064	Y

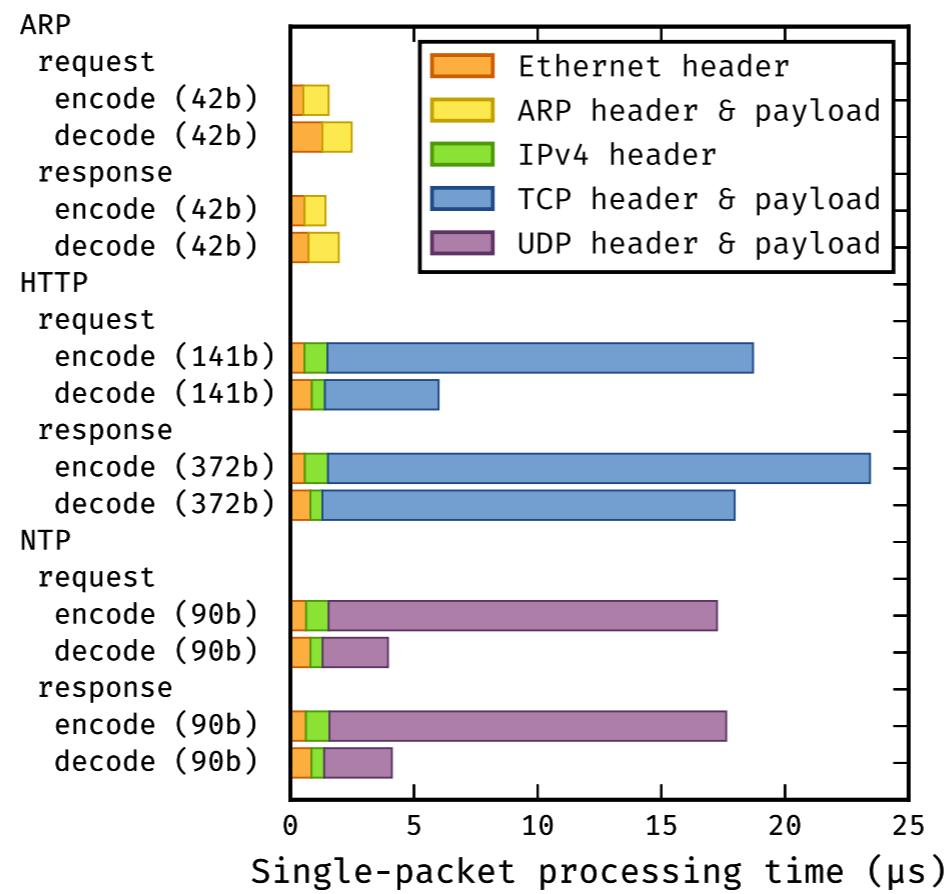
Component Library

# Narcissus in Action

- MirageOS is a library operating system for secure, high-performance network applications written in OCaml
- Replaced IP stack of MirageOS with extracted OCaml implementations of synthesized decoders.

Protocol	LoC	Interesting Features
Ethernet	150	Multiple format versions
ARP	41	
IP	141	IP Checksum; underspecified fields
UDP	115	IP Checksum with pseudoheader
TCP	181	IP Checksum with pseudoheader; underspecified fields
DNS	474	DNS compression; variant types

Derived Decoders



# Shameless Plug

- PL @Purdue:
  - 10 faculty across 2 departments
  - Many interests Formal Methods, Synthesis, Compilers, Metaprogramming, PL Theory, Static Analysis, PL+Security, PL+ML ...



- Undergrads+Master's students: Please applying!
- Follow us on twitter + facebook:
  - Twitter: @purdue\_pl
  - <https://www.facebook.com/PurduePL/>

# Conclusion

- Coq is a powerful tool for **automatically** deriving correct software
  - Powerful specification language (formats + correctness)
  - Gallina is a pretty good functional language (implementations)
  - Tactics are useful for automating construction of well-understood terms (derivations)
- If you're working on verified distributed systems; let's talk!  
(Happy to talk if you're not working on distributed systems too)

Questions?

