# From Verified Parsers and Serializers to Format-Aware Fuzzers

Benjamin Delaware
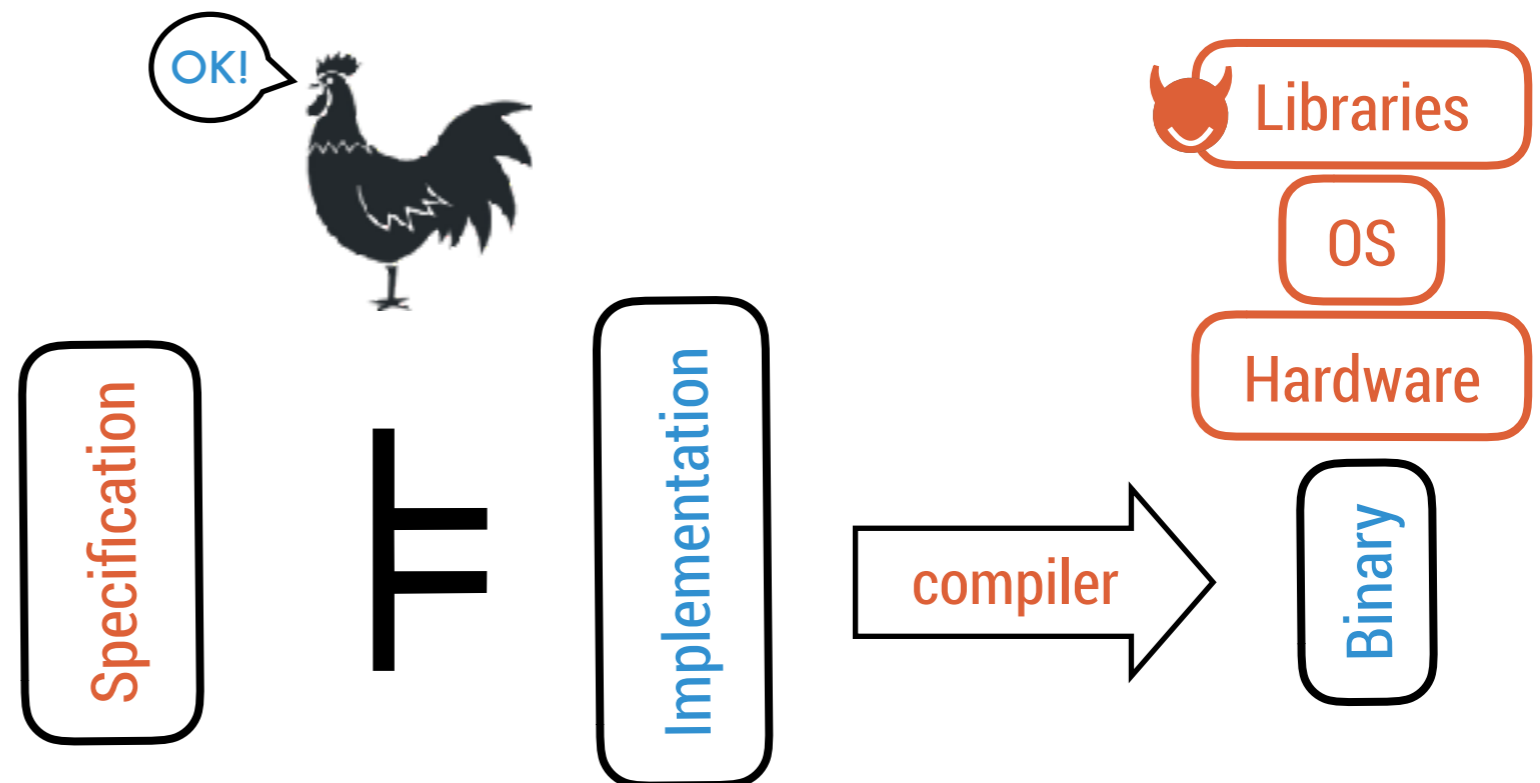
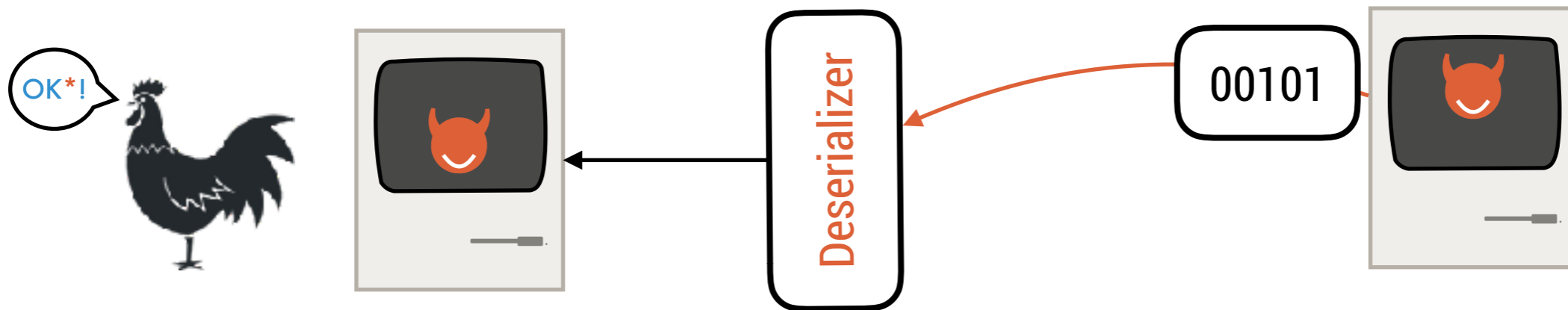Purdue Computer Science

# Formal Verification

- Numerous developments of high-assurance software in proof assistants in the past five years:
  - CompCert C compiler
  - seL4 microkernel
  - FSCQ file system

- Assurance comes from formal guarantees[*] provided by proof assistant:
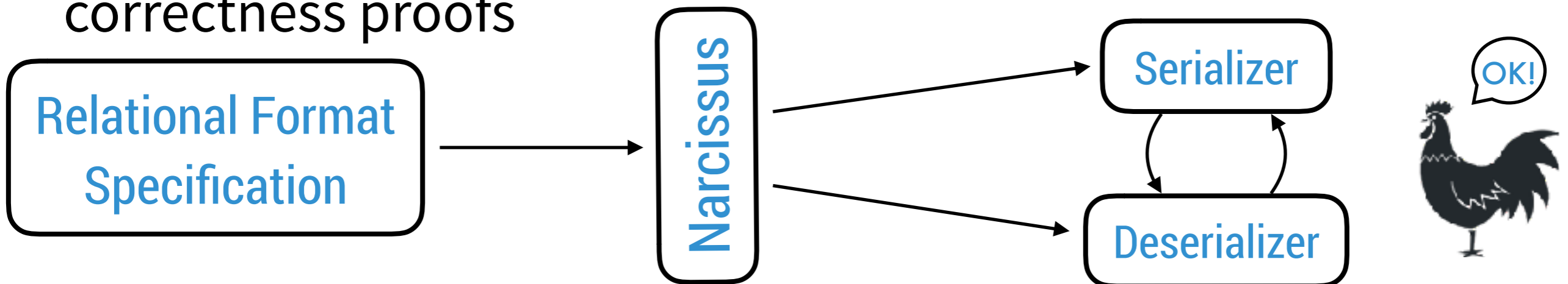
[*] w.r.t Trusted Base

# Narcissus

- For networked systems, deserialization is important[1]
  - If these are in your TCB, bugs will break the assurance case!



- Enter Narcissus:

  - User-extensible framework for synthesizing encoders and decoders from format specifications, with machine-checked correctness proofs
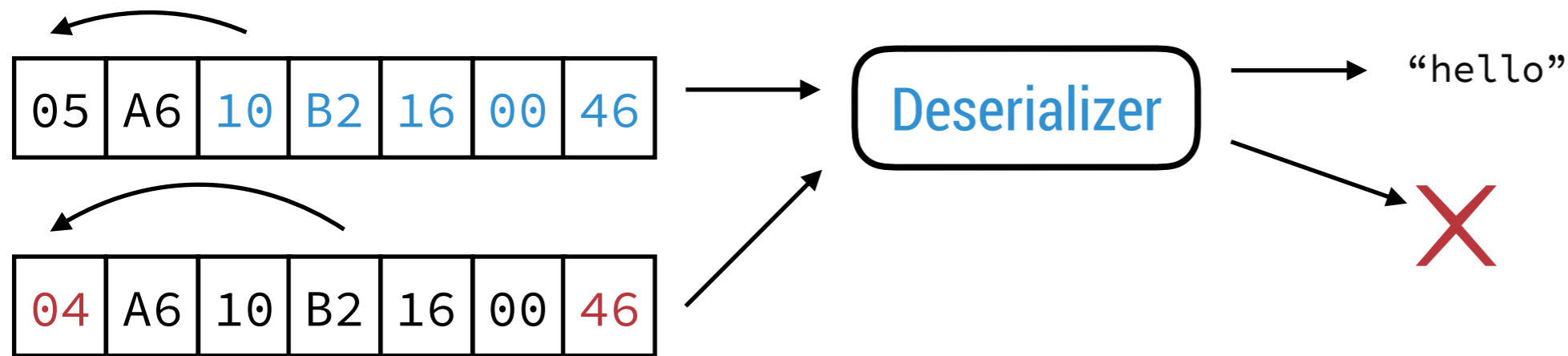


[1] An Empirical Study on the Correctness of Formally Verified Distributed Systems. Pedro Fonseca, Kaiyuan Zhang, Xi Wang, and Arvind Krishnamurthy.

# All Done?

- *Probably* unreasonable to incorporate synthesized decoders and decoders into every existing codebase.
  - Synthesized code is OCaml (working on verified C)
  - Assumes clean interface between communication and processing code
- How to leverage work to secure legacy code?

# From Verification to Fuzzing

- Formats can contain implicit dependencies
- These decoders are provably correct recognizers for the *entire* input format.



- Verification exposes latent dependencies in formats.
- Hypothesis: these dependencies can be leveraged to generate format-aware fuzzers.
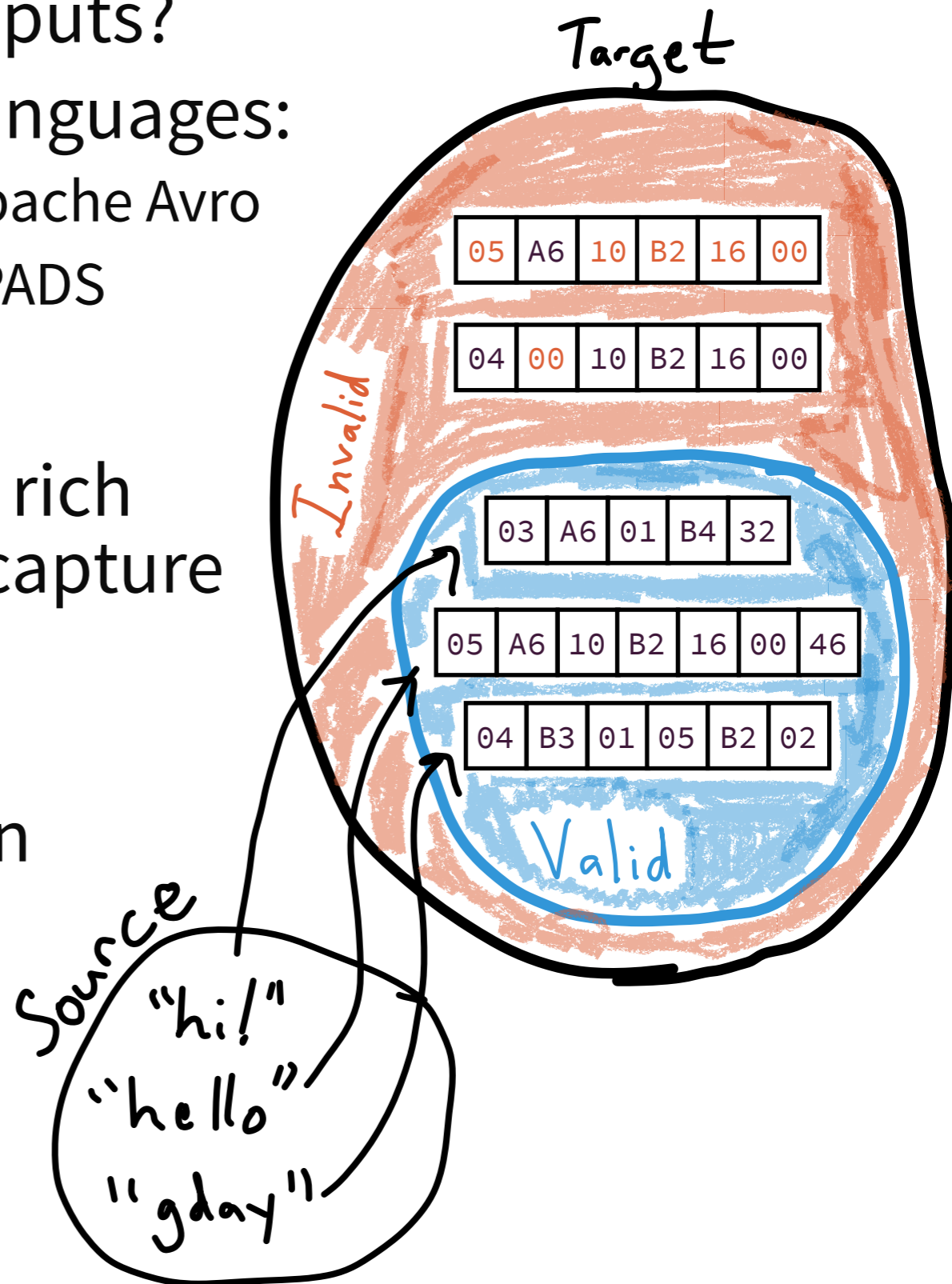
# Today's Talk

- Embedding Formats in Narcissus
- Synthesizing Correct-by-Construction encoders and decoders
- Leveraging these to generate format-aware fuzzers

# Specifying Formats in Narcissus

- **First challenge**: specifying valid inputs?
- Established format specification languages:
  - Interface Generators: ASN.1, Protobuffs, Apache Avro
  - Format Specification Languages: binpac, PADS

- Internet servers were the original verification target, so we needed a rich enough specification language to capture legacy formats.

- **Solution (?):** functional description
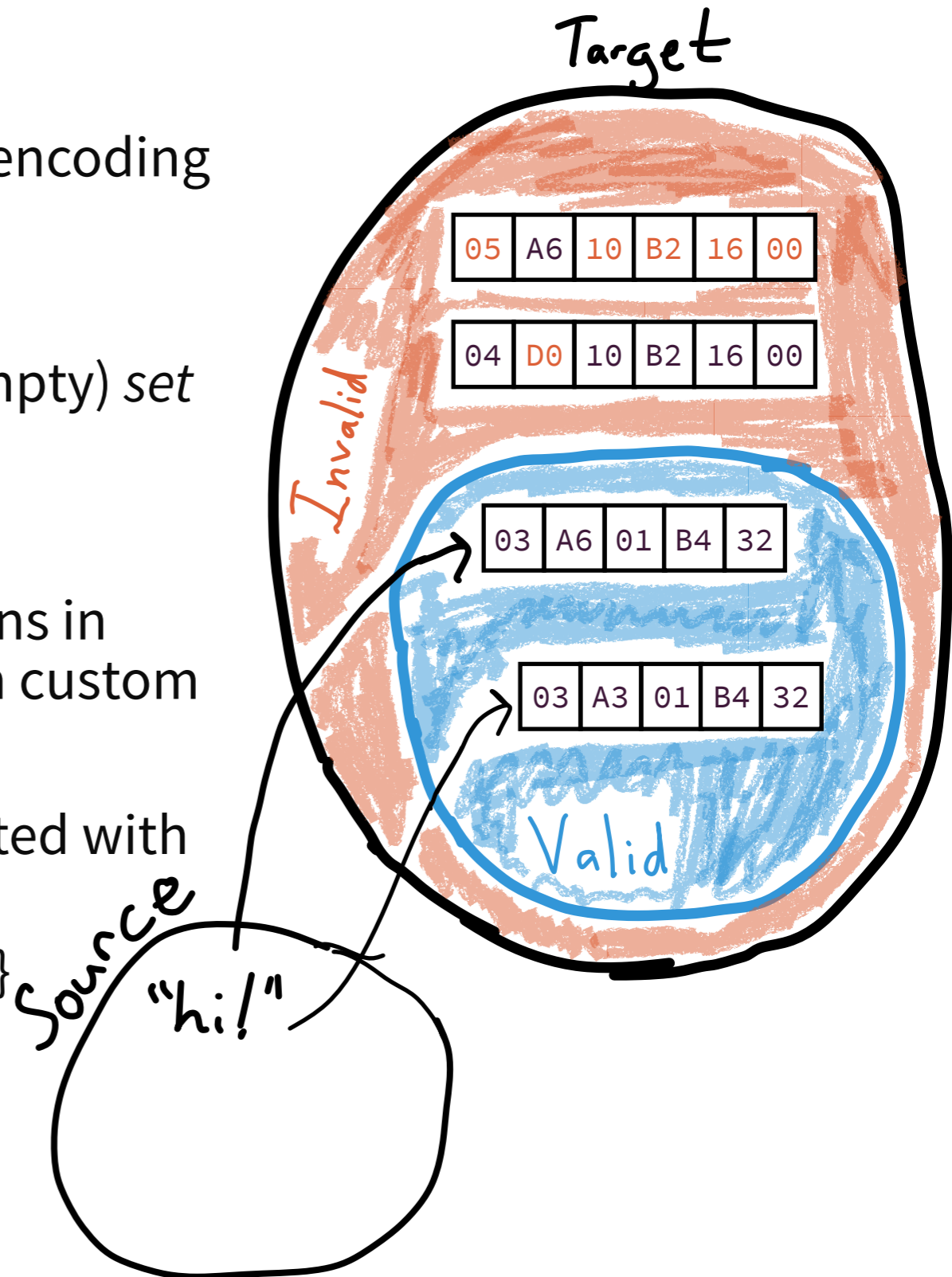
```
format(s) = |s| ++ 166 ++ s
```

# Relational Specifications

- Many formats do not have a single canonical encoding of a source value
  - i.e. DNS packet compression

- **Solution**: map source values to a (possibly empty) *set* of target representations:

$$\text{format}(s) = |s| + \{n \mid n \le 2^{17}\} + s$$

- These relations are represented as propositions in Coq's logic, so users can freely write their own custom format specifications

- Constraints on source values can be represented with set intersection:

$$\text{format'}(s) = \text{format}(s) \cap \{(s,t) \mid |s| \le 2^{17}\}$$

Target

| 05 | A6 | 10 | B2 | 16 | 00 |

| 04 | D0 | 10 | B2 | 16 | 00 |

Invalid

| 03 | A6 | 01 | B4 | 32 |

| 03 | A3 | 01 | B4 | 32 |

Valid

Source

"hi!"

# Simplifying Specifications

- Narcissus includes a library of common formats
  - <u>Base formats</u> for single data types
  - <u>Combinators</u> for composing formats

| 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

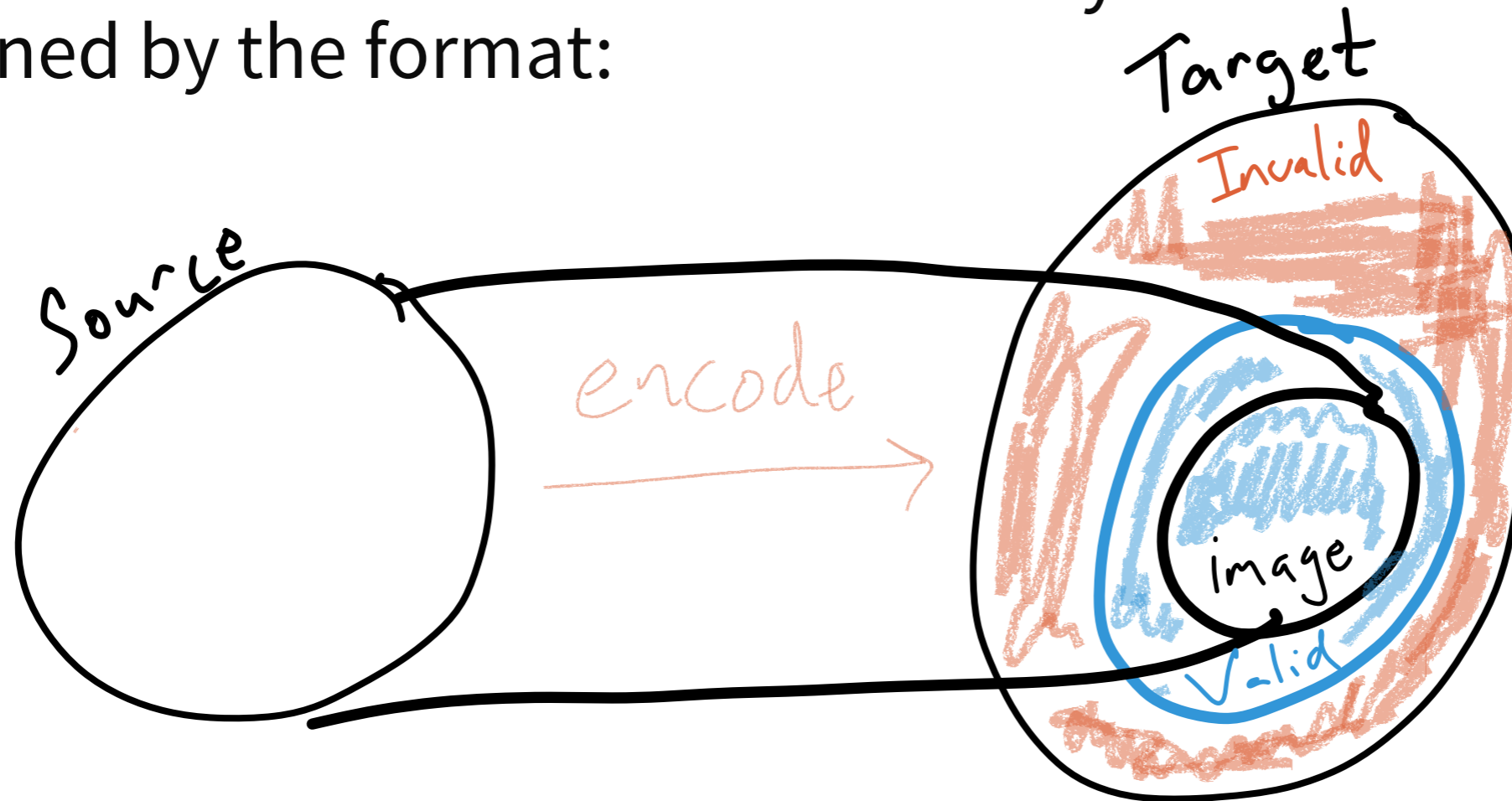# Simplifying Specifications

- Narcissus includes a library of common formats
  - <u>Base formats</u> for single data types
  - <u>Combinators</u> for composing formats

```
Definition IPv4_Packet_Format (ip4 : IPv4_Packet) :=
    format_nat 4 4 # format_nat 4 (5 + |ip4.Options|) # {n : char | true}
  # format_word ip4.TotalLength
  # format_word ip4.ID
  # {b : bool | true} # format_bool ip4.DF # format_bool ip4.MF # format_word ip4.FragmentOffset
  # format_word ip4.TTL # format_enum ProtocolCodes ip4.Protocol
  # IPChecksum_Valid
  # format_word ip4.SourceAddress
  # format_word ip4.DestAddress
  # format_list format_word ip4.Options # e.
```

| Bits | 0–3 | 4–7 | 8–11 | 12–15 | 16–18 | 19–23 | 24–27 | 28–31 |
|------|-----|-----|------|-------|-------|-------|-------|-------|
| 0 | Version | Head Length | Type of Service | | Total length (Packet) | | | |
| 32 | Identification | | | | Flags | Fragment Spacing | | |
| 64 | Lifespan (TTL) | | Protocol | | Header checksum | | | |
| 96 | Source Address | | | | | | | |
| 128 | Destination Address | | | | | | | |
| 160 | Options | | | | | | | |

# Specifying Encoders and Decoders

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



```
EncoderOK(Format, e) ≡ ∀s.Format ∋ (s, e(s))
```

# Specifying Encoders and 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
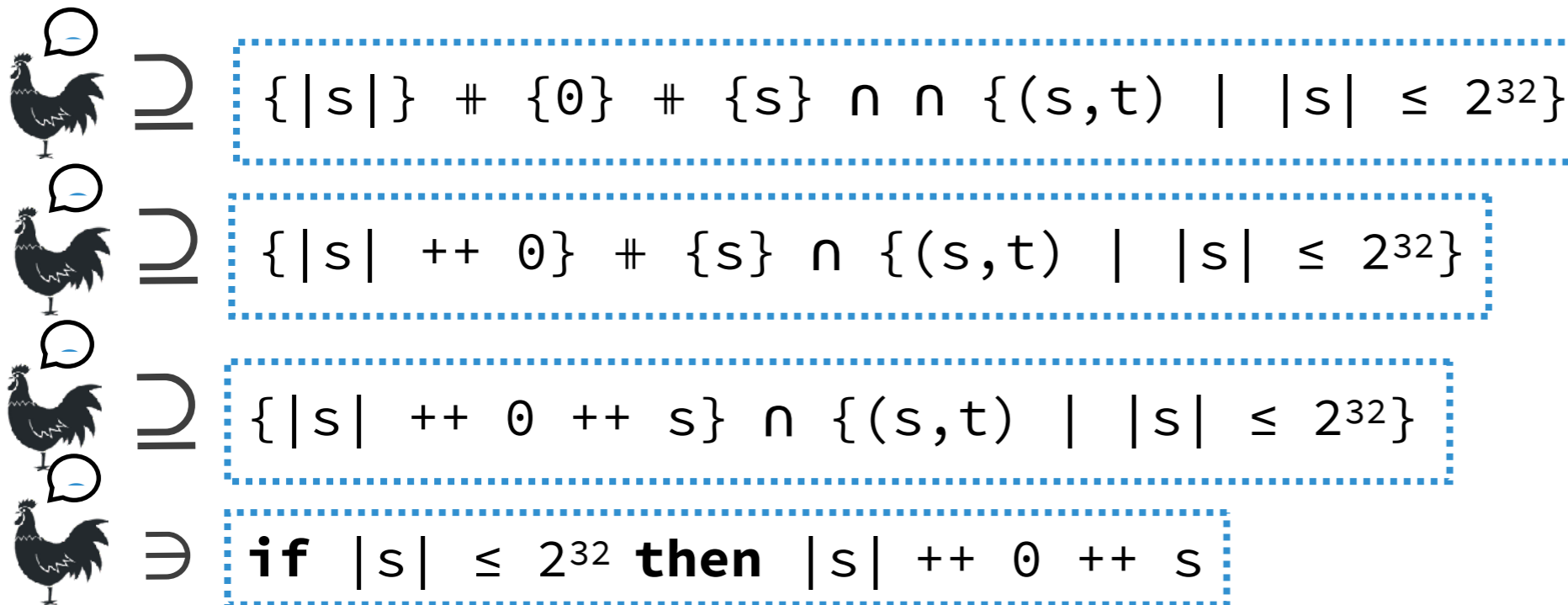


```
DecoderOK(Format, d) ≡ ∀t.Format ∋ (d(t), t)
                       ∧ d(t) = ⊥ → ∀v. Format ∌ (v, t)
```

# Deriving Encoders

- Can phrase construction of a correct encoder as a user directed search for a function satisfying `EncoderOK`
  - Such searches are the bread and butter of theorem provers

- <u>Key Observation</u>: formats are inherently compositional, so this process can be decomposed into a series of small steps
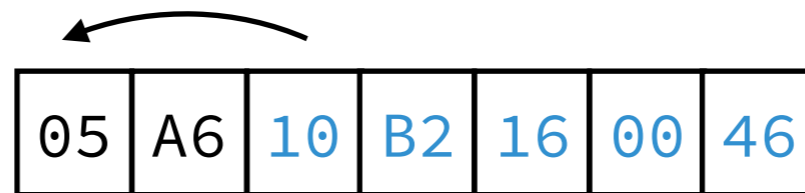
```
format'(s) := {|s|} # {n | n ≤ 2¹⁷} # {s} ∩ {(s,t) | |s| ≤ 2³²}
```

$$\supseteq \quad \{|s|\} \# \{0\} \# \{s\} \cap \cap \{(s,t) \mid |s| \le 2^{32}\}$$

$$\supseteq \quad \{|s| \text{ ++ } 0\} \# \{s\} \cap \{(s,t) \mid |s| \le 2^{32}\}$$

$$\supseteq \quad \{|s| \text{ ++ } 0 \text{ ++ } s\} \cap \{(s,t) \mid |s| \le 2^{32}\}$$

$$\ni \quad \textbf{if } |s| \le 2^{32} \textbf{ then } |s| \text{ ++ } 0 \text{ ++ } s$$

- These proofs can be automated

# Deriving Decoders

- Can do the same for decoders, but correctness of subdecoders now depends on other parts of the encoded value:



| 05 | A6 | 10 | B2 | 16 | 00 | 46 |

- DNS— compressed domains are pointers
- DNS— resource record tag determines how payload is parsed
- SDN— versions effects available options
- ZIP— position of start of central directory depends on EOCD

```
∀n. DecoderOK({s} ∩ {(s,t) | |s| = n}, decodeList n)

where decode 0 [] = Some []
      decode n (c : t) = decode (n − 1) t >>= \l -> c : l
      decode _ _   = None
```

# Deriving Decoders[2]

- <u>Key idea</u>: keep track of dependence data when decomposing proof:



```
  DecoderOK(Format₁', d₁)
∧ image(Format₁') = image(Format₁)
∧ DecoderOK(Format₂ ∩ {(s,t) | ∃t'. (v, t') ∈ Format₁'
                                ∧       (s, t') ∈ Format₁},
       d₂(v) )
 → DecoderOK(Format₁ # Format₂, d₁ >>= d₂)
```

# Deriving Decoders[2]

- <u>Key idea</u>: keep track of dependence data when decomposing proof:

DecoderOK({|s|} ⊎ {n | n ≤ $2^{17}$} ⊎ {s} ∩ {(s,t) | |s| ≤ $2^{32}$}, ?)

DecoderOK({n | n ≤ $2^{17}$} ⊎ {s} ∩ {(s,t) | |s| ≤ $2^{32}$} ∩ {v = |s|}, ? v)

DecoderOK({s} ∩ {(s,t) | |s| ≤ $2^{32}$} ∩ {v = s} ∩ {n ≤ $2^{17}$}, ? v n)

DecoderOK({(s,t) | |s| ≤ $2^{32}$} ∩ {v = |s|} ∩ {n ≤ $2^{17}$} ∩ {l = s}, ? v n l)

DecoderOK({(s,t) | |s| ≤ $2^{32}$ ∧ v = |s| s ∧ ≤ $2^{17}$ ∧ l = s}, l)

# Deriving Decoders[2]

- <u>Key idea</u>: keep track of dependence data when decomposing proof:

```
DecoderOK({|s|} ⊎ {n | n ≤ 2¹⁷} ⊎ {s} ∩ {(s,t) | |s| ≤ 2³²},
            v <- decodeChar;
            n <- decodeChar;
            l <- decodeList v;
            if n <= 2¹⁷ then return l else None)
```

# Narcissus in Action

- MirageOS is a library operating system for secure, high-performance network applications written in OCaml

- Replaced network stack of MirageOS with extracted OCaml implementations of synthesized decoders.

- Found one problem in the test suite.

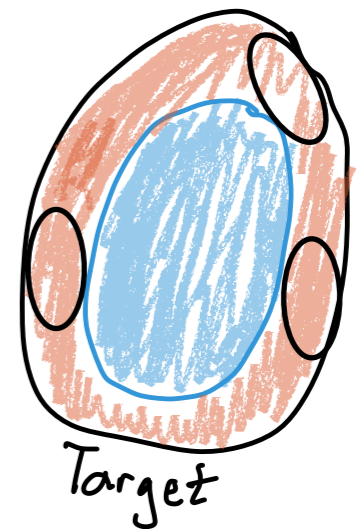| Protocol | LoC | Interesting Features |
|----------|-----|----------------------|
| Ethernet | 150 | Multiple format versions |
| ARP | 41 | |
| IP | 141 | IP Checksum; underspecified fields |
| UDP | 115 | IP Checksum with pseudoheader |
| TCP | 181 | IP Checksum with pseudoheader; under-specified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

Derived Decoders

- But, *probably* unreasonable to incorporate synthesized decoders and decoders into every existing codebase.

- How can we leverage this to secure legacy systems?

# Towards Format-Aware Fuzzers

- The final decoder synthesis step contains the accumulated dependencies embedded in the format:

  `DecoderOK({(s,t) | |s| ≤ 2³² ∧ n ≤ 2¹⁷ ∧ v = |s| ∧ l = s}, ?)`

  - invariants on the original input data
  - invariants on the shape of the target values
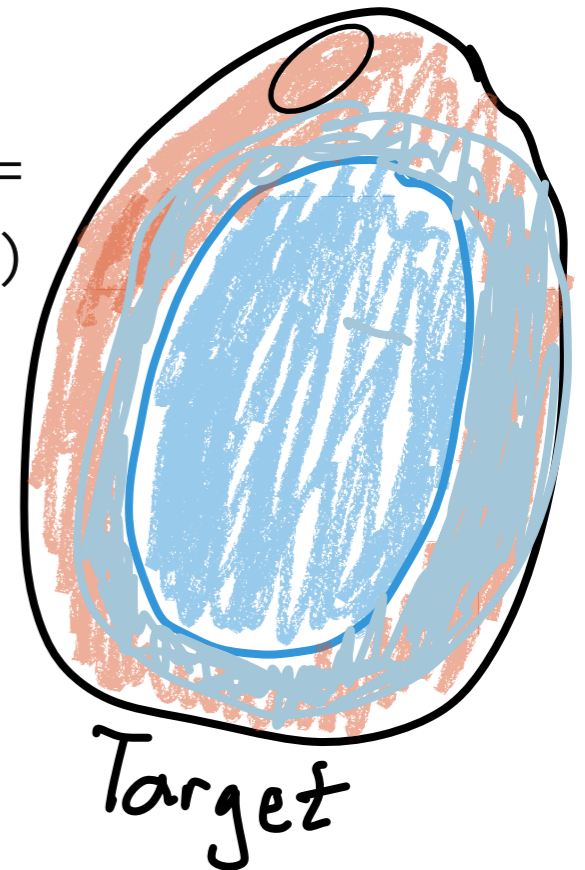  - dependencies between bytes of the target values



Target

- <u>Idea</u>: violating any one of these these dependencies yields an input not included in the format
- Can we selectively break these dependencies to "fuzz" the format in a smart way?
- Generate predicates for behavioral property testing?

# Gradual Fuzzing

- We don't need to formalize the full format to get useful fuzzers:

  - Only specifying certain fields tests dependencies between these fields

  - Rest of the target value is "don't care" bits:

```
Definition IPv4_Packet_Format (ip4 : IPv4_Packet) :=
    format_nat 4 4 # format_nat 4 (5 + |ip4.Options|)
 # {n : char | true}
 # {n : 16 words | true}
 # format_list format_word ip4.Options # e.
```

- Gradually specify complex formats, hitting low-hanging bits first
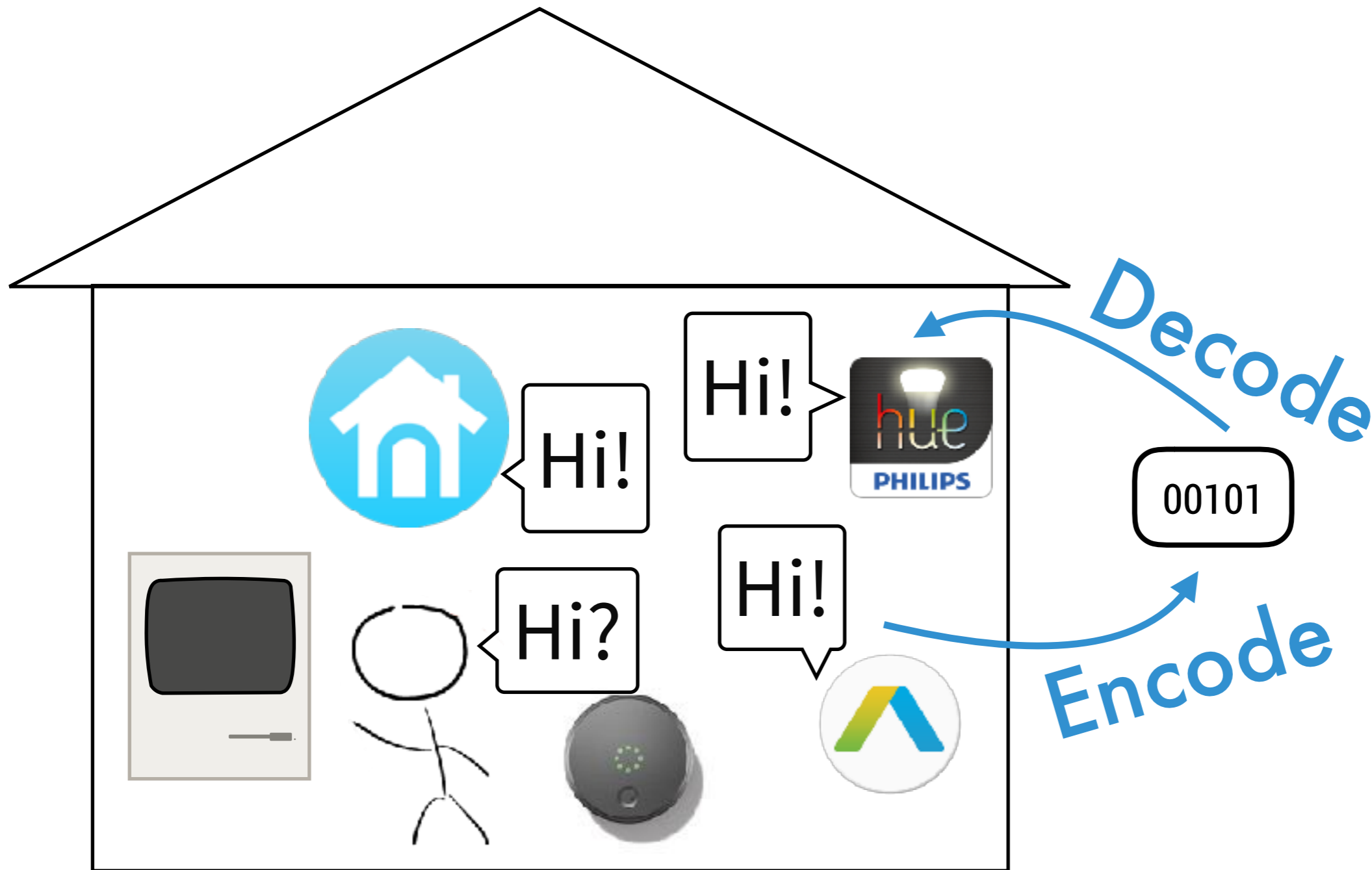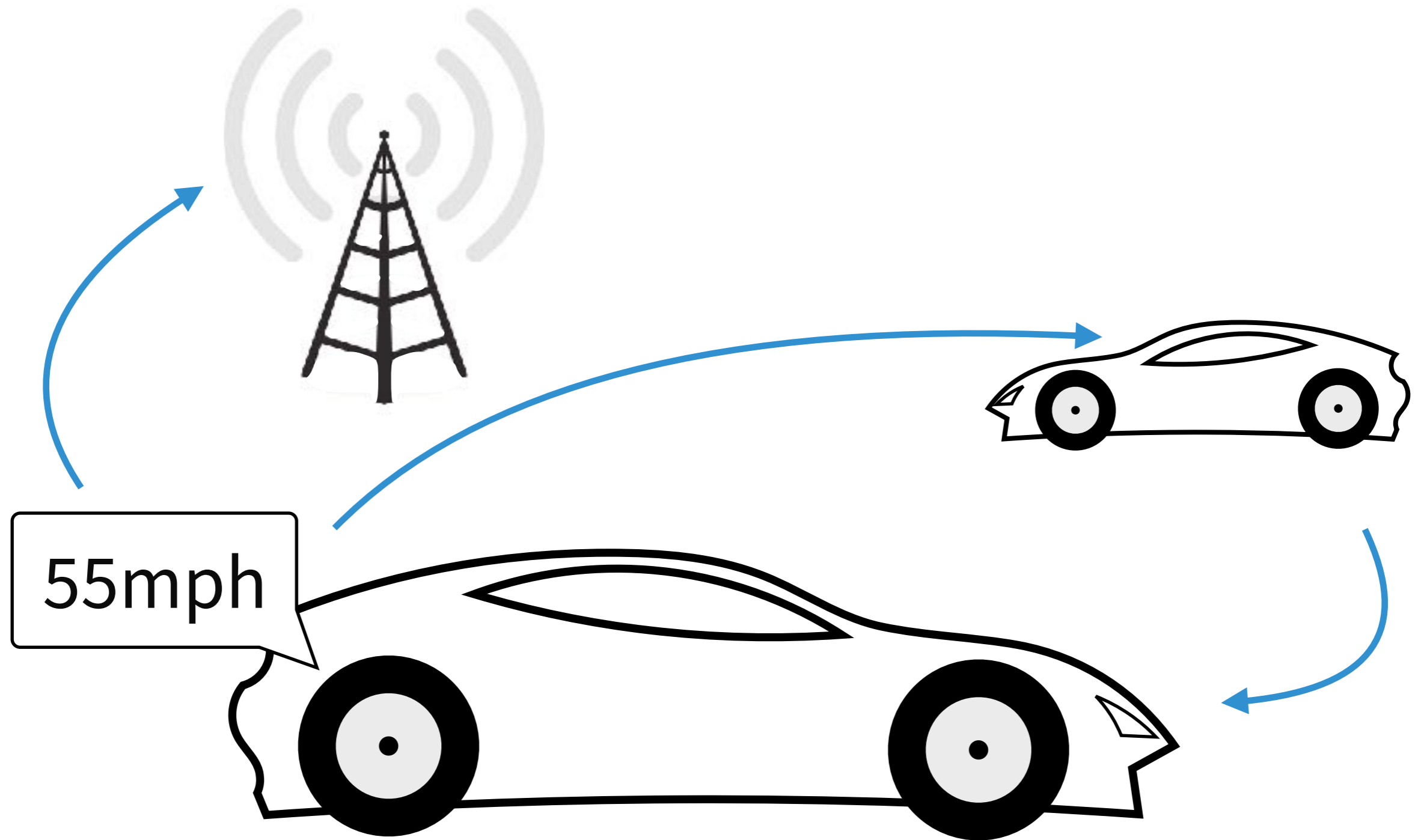


Target

# Conclusion

- Today's talk:
  - Embedding Formats in Narcissus
  - Synthesizing Correct-by-Construction encoders and decoders
  - Leveraging these to generate format-aware fuzzers

## Thoughts?

# Conclusion

- Today's talk:
  - Embedding Formats in Narcissus
  - Synthesizing Correct-by-Construction encoders and decoders
  - Leveraging these to generate format-aware fuzzers

- Next Steps:
  - Evaluation?
  - Thoughts?

# Conclusion

- Today's talk:
  - Embedding Formats in Narcissus
  - Synthesizing Correct-by-Construction encoders and decoders
  - Leveraging these to generate format-aware fuzzers

- Next Steps:
  - Evaluation?
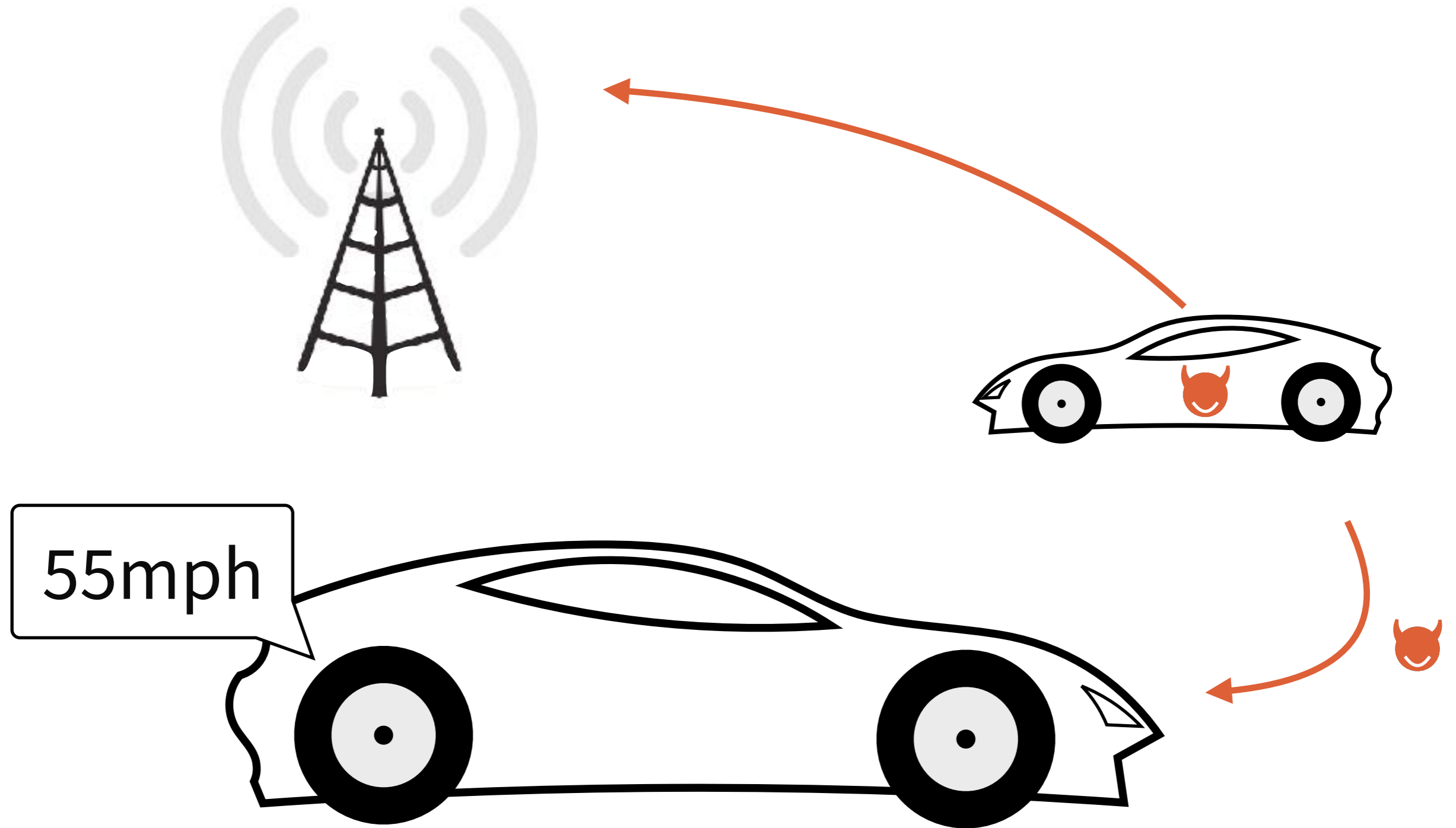  - Thoughts?

# Computers are Multiplying



- Relationship between encoded + decoded data important
- Bugs lead to miscommunication

# Communication is Multiplying



55mph

- Decoders present attack surface for malicious packets
- NTSB will likely mandate V2V communication within the next decade

# Communication is Multiplying



55mph

- Decoders present attack surface for malicious packets
- NTSB will likely mandate V2V communication within the next decade
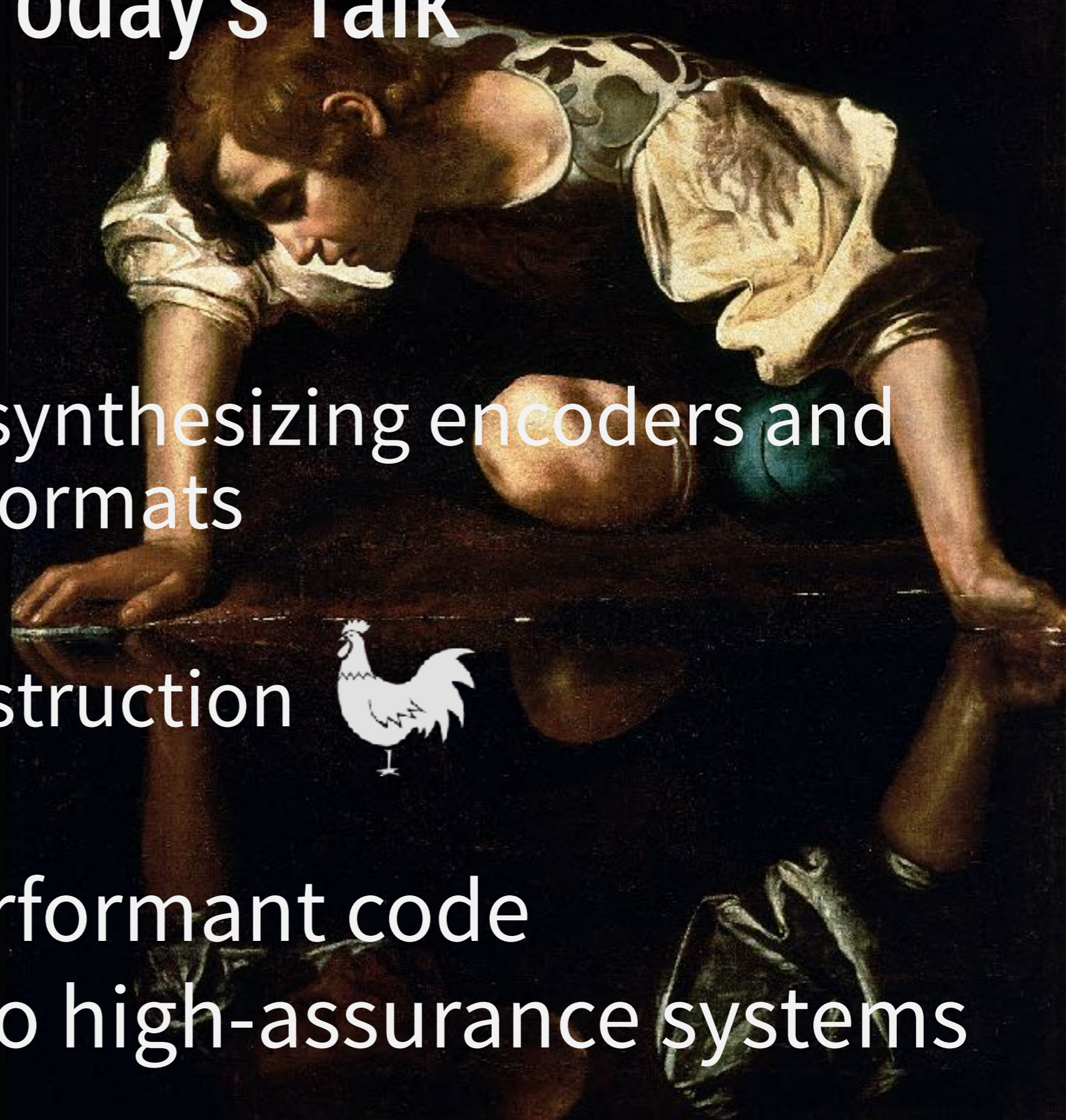
# Why Worry?

Since 2013:

And Many More!

# Established Solutions

- Interface Generators:
  - ASN.1, Protobuffs, Apache Avro
  - Data format defined by system

- Format Specification Languages:
  - binpac, PADS, Packet Types
  - New formats still require modifying code generator

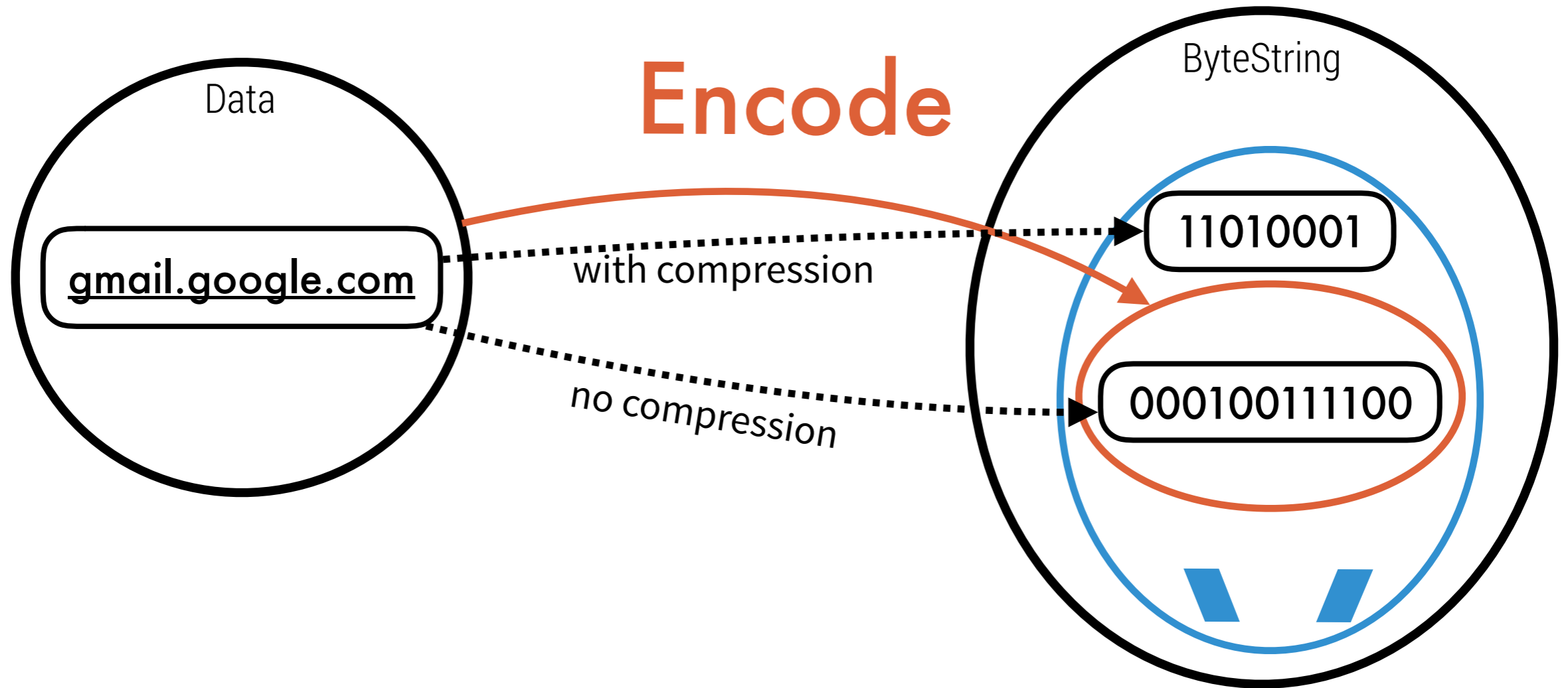- User-Extensible Systems:
  - Nail
  - No formal guarantees

# Today's Talk

- ## Narcissus:
  - Framework for synthesizing encoders and decoders from formats
  - User extensible
  - Correct-by-Construction

- Generating performant code
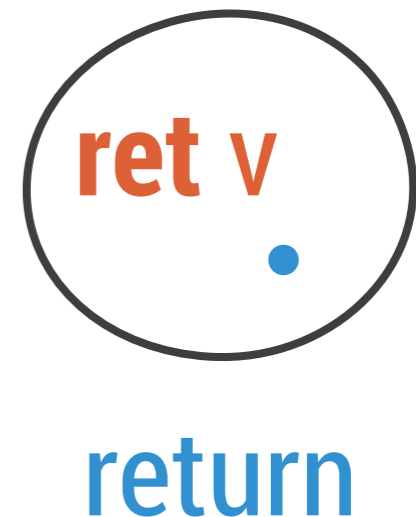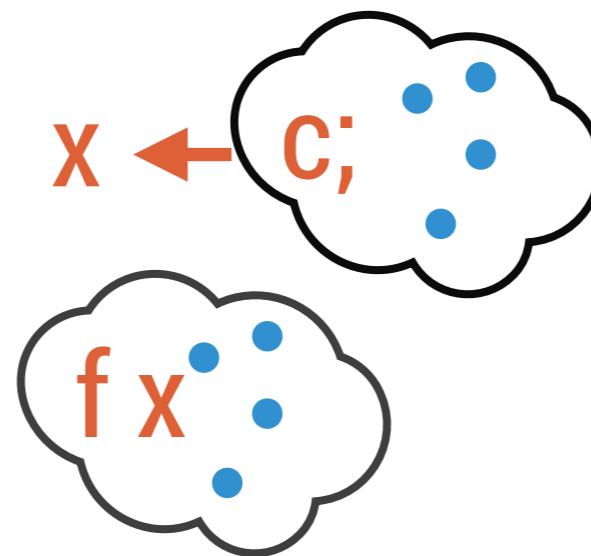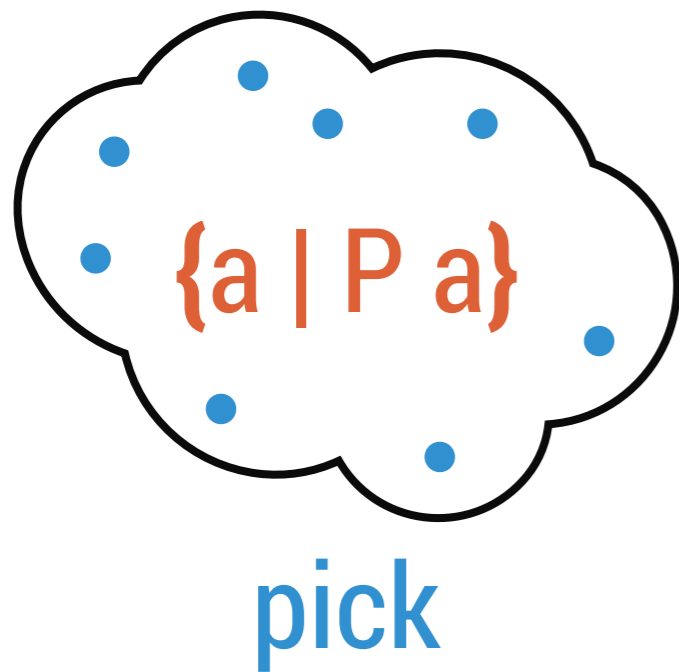- Integration into high-assurance systems

# Specifying Formats



Data

Encode

ByteString

gmail.google.com

with compression

11010001

no compression

000100111100

- Mapping is not one-to-one: compression, unspecified fields
- Key Idea: specify set of valid encodings for value as a binary relation
- Encoder always maps into valid set

# Specifying Formats

Key Idea: Represent formats as functional programs in the nondeterminism monad.



pick

bind

return

# Computations

> **Key Idea:** Represent formats as functional programs in the nondeterminism monad.

```
 0  1  2  3  4  5  6  7  8
+--+--+--+--+--+--+--+--+--+
|          NUMREADINGS     |
+--+--+--+--+--+--+--+--+--|
|                          |
/            ID            /
|                          |
+--+--+--+--+--+--+--+--+--+
|    CLASS        |0 |0 |0 |
+--+--+--+--+--+--+--+--+--+
|                          |
/          READINGS        /
|                          |
+--+--+--+--+--+--+--+--+--+
```
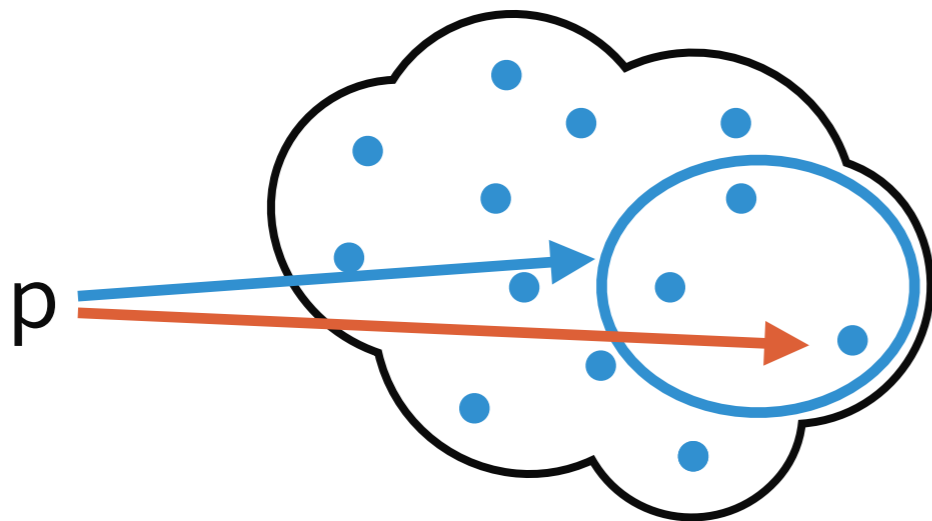
```
Packet := ⟨ID :: string,
            readings :: list word⟩

SimpleFormat (p : Packet) :=
 b₁ ← formatNat |p!readings|;
 b₂ ← formatString p!ID;
 b₃ ← {w : word | w < 32};
 b₄ ← formatList encodeWord p!readings;
 ret (b₁⧺b₂⧺b₃⧺b₄)
```

# Specifying Correct Encoders

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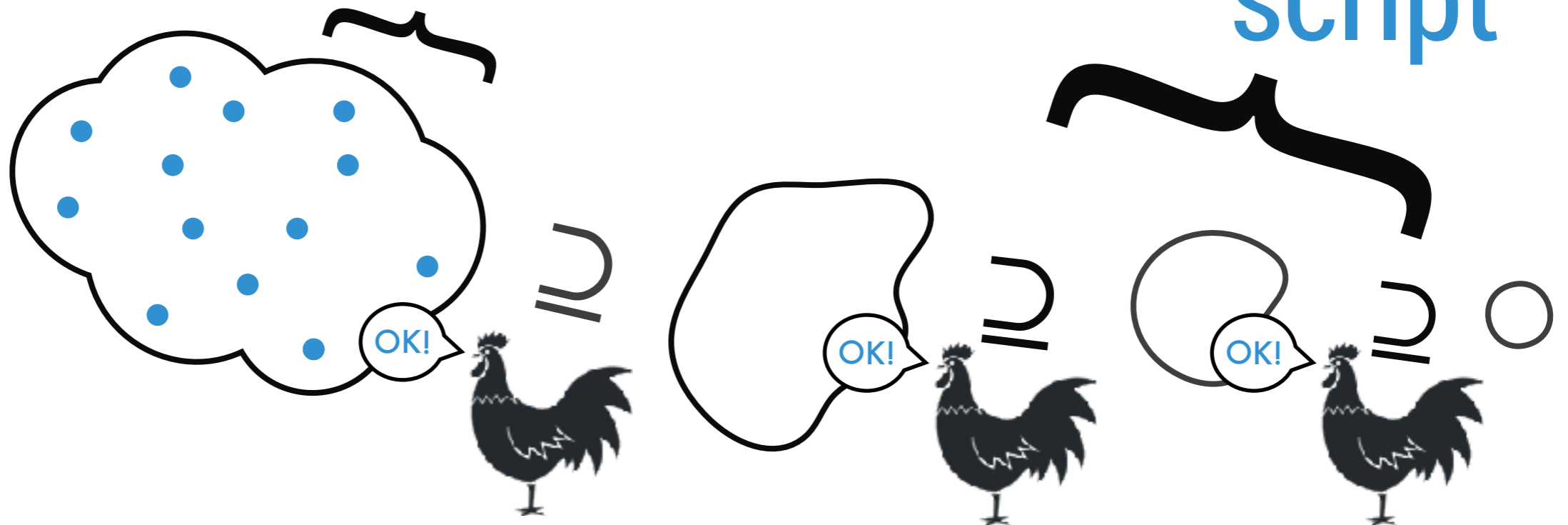


∀p. SimpleFormat(p) ⊇ SimpleEncoder(p)

# Deriving Correct Encoders

The construction of a correct encoder can be posed as a user-guided search in a proof assistant.

format

optimization script

# Properties of Refinement

- Preorder

$$\frac{a \supseteq b \quad b \supseteq c}{a \supseteq c} \text{ Trans}_{\supseteq} \qquad\qquad \frac{}{a \supseteq a} \text{ Refl}_{\supseteq}$$

- Respected by sequencing

$$\frac{a \supseteq b}{r \leftarrow a; f(r) \supseteq r \leftarrow b; f(r)} \text{Seq1}_{\supseteq}$$

$$\frac{\forall r, f(r) \supseteq f'(r)}{r \leftarrow a; f(r) \supseteq r \leftarrow a; f'(r)} \text{Seq2}_{\supseteq}$$

# Deriving Correct Encoders



rewrite
formatNatOK!

```
SimpleFormat (p : Packet) :=
  b₁ ← formatNat |p!readings|;
  b₂ ← formatString p!ID;
  b₃ ← {w : word | w < 32};
  b₄ ← formatList encodeWord p!readings;
  ret (b₁⧺b₂⧺b₃⧺b₄)
```

∣∪

```
SimpleFormat (p : Packet) :=
  b₁ ← encodeNat |p!readings|;
  b₂ ← formatString p!ID;
  b₃ ← {w : word | w < 32};
  b₄ ← formatList encodeWord p!readings;
  ret (b₁⧺b₂⧺b₃⧺b₄)
```

# Deriving Correct Encoders



rewrite
formatStrOK!

```
SimpleFormat (p : Packet) :=
 b₁ ← encodeNat |p!readings|;
 b₂ ← formatString p!ID;
 b₃ ← {w : word | w < 32};
 b₄ ← formatList encodeWord p!readings;
 ret (b₁#b₂#b₃#b₄)
```

∐

```
SimpleFormat (p : Packet) :=
 b₁ ← encodeNat |p!readings|;
 b₂ ← encodeString p!ID;
 b₃ ← {w : word | w < 32};
 b₄ ← formatList encodeWord p!readings;
 ret (b₁#b₂#b₃#b₄)
```

# Deriving Correct Encoders



rewrite
MyRule!

```
SimpleFormat (p : Packet) :=
  b₁ ← encodeNat |p!readings|;
  b₂ ← encodeString p!ID;
  b₃ ← {w : word | w < 32};
  b₄ ← formatList encodeWord p!readings;
  ret (b₁#b₂#b₃#b₄)
```

$\cup$

```
SimpleFormat (p : Packet) :=
  b₁ ← encodeNat |p!readings|;
  b₂ ← encodeString p!ID;
  b₃ ← ret 0;
  b₄ ← formatList encodeWord p!readings;
  ret (b₁#b₂#b₃#b₄)
```
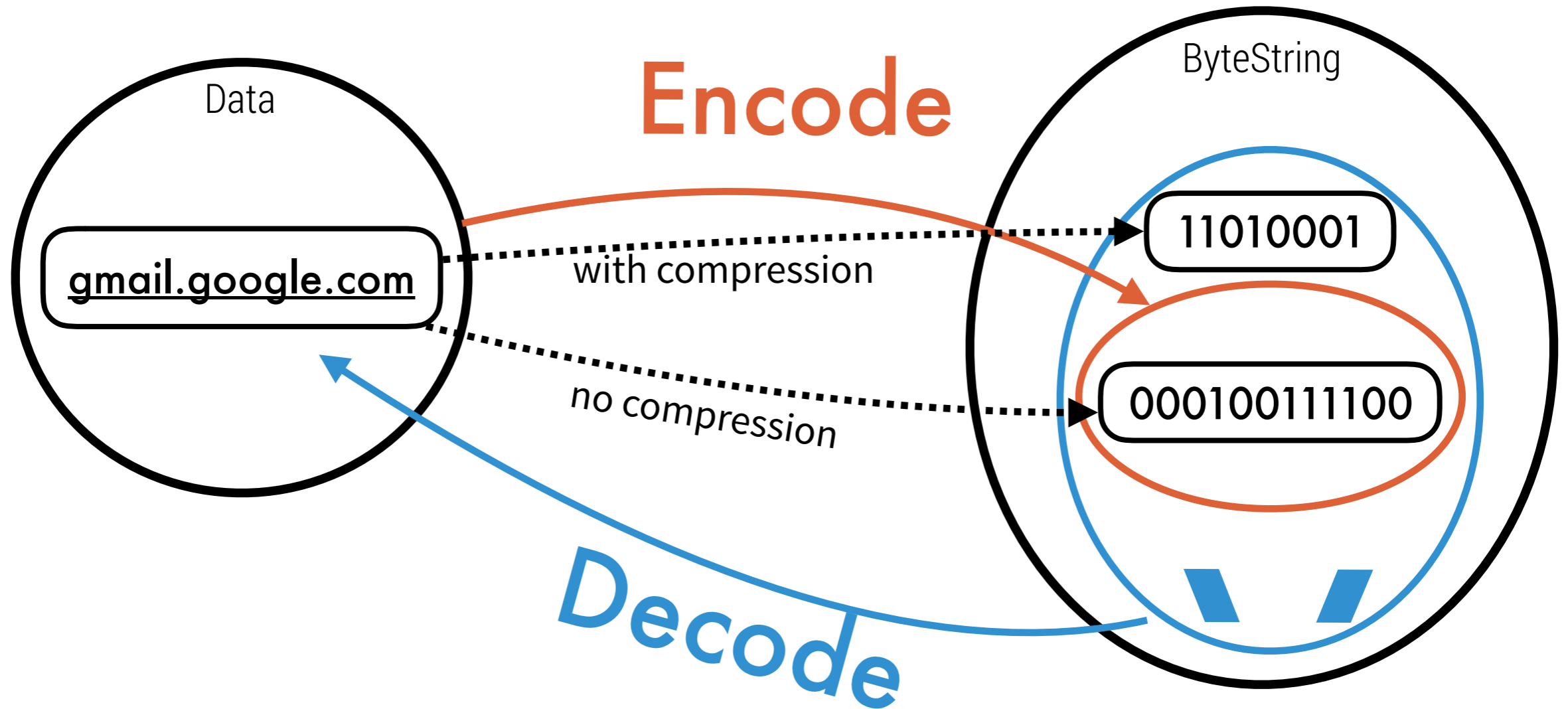
# Deriving Correct Encoders

finish!

```
SimpleFormat (p : Packet) :=
  b₁ ← encodeNat |p!readings|;
  b₂ ← encodeString p!ID;
  b₃ ← ret 0;
  b₄ ← formatList encodeWord p!readings;
  ret (b₁#b₂#b₃#b₄)
```

⊎

```
SimpleEncoder (p : Packet) :=
encodeNat |p!readings| # encodeString p!ID # 0
# encodeList encodeWord p!readings
```

- Users can safely add their own formats and rewrite rules
- Rewrites can be packaged together into single optimization tactic

# Specifying Correct Decoders



Data

gmail.google.com

Encode

with compression

no compression

ByteString

11010001

000100111100

Decode

$\text{Valid}^{-1}\, b\, P \triangleq \{\, p \mid b \in \text{Valid}\, p \wedge P\, p \; \wedge \; \neg\exists\, p.\, b \in \text{Valid}\, p \wedge P\, p \rightarrow p = \bot\, \}$

# Deriving Correct Decoders

The construction of a correct decoder can **also** be posed as a user-guided search in a proof assistant.



$$\frac{}{\texttt{formatNat}^{-1}\ \texttt{b}\ \mathbb{T}\ \ni\ \texttt{decodeNat(b)}}$$

Component Library

$$\frac{}{\texttt{formatString}^{-1}\texttt{(b)}\ \mathbb{T}\ \ni}$$

Invariant on List Elements

Know Length

$$\frac{\texttt{format}_A^{-1}\ \texttt{b}\ P_A\ \ni\ \texttt{decode}_A\texttt{(b)}\quad \mathbf{Q}(l)\ \rightarrow\ \forall a \in l.\ P_A(a)\qquad \mathbf{Q}(l)\rightarrow|l|\ =\ n}{\texttt{formatList}^{-1}\ \texttt{format}_A\ \texttt{b}\ \mathbf{Q}\ \ni\ \texttt{decodeList decode}_A\texttt{(b, n)}}$$

# Deriving Correct Decoders

The construction of a correct decoder can **also** be posed as a user-guided search in a proof assistant.

Component Library

$\forall a.\ P_A(a) \rightarrow format_B^{-1}\ b\ Q \ni decode_B(a,\ b)$

$format_A^{-1}\ b\ P_A \ni decode_A(b)$ $\qquad\qquad$ $Q(ab) \rightarrow P_A(\pi\ ab)$

─────────────────────────────────────────────

$format_A;\ format_B\ ^{-1}\ b\ Q \ni (b',a) \leftarrow decode_A\ b\ P_A;$

$decode_B\ (a,b')\ Q$

$\forall a'.\ Q(a') \rightarrow a'=a$

─────────────────────────────

**ret** $[]^{-1}\ b\ Q \ni Some\ a$

# Deriving Correct Decoders

rewrite
DecodeStrOK!

```
SimpleDecoder (b : ByteString) :=
SimpleFormat⁻¹ b 𝕋
```

⊔

```
SimpleDecoder (b : ByteString) :=
 (n, b) ← decodeNat(b); ???
```

⊔

```
SimpleDecoder (b : ByteString) :=
 (n, b) ← decodeNat(b);
 (s, b) ← decodeString(b); ???
```

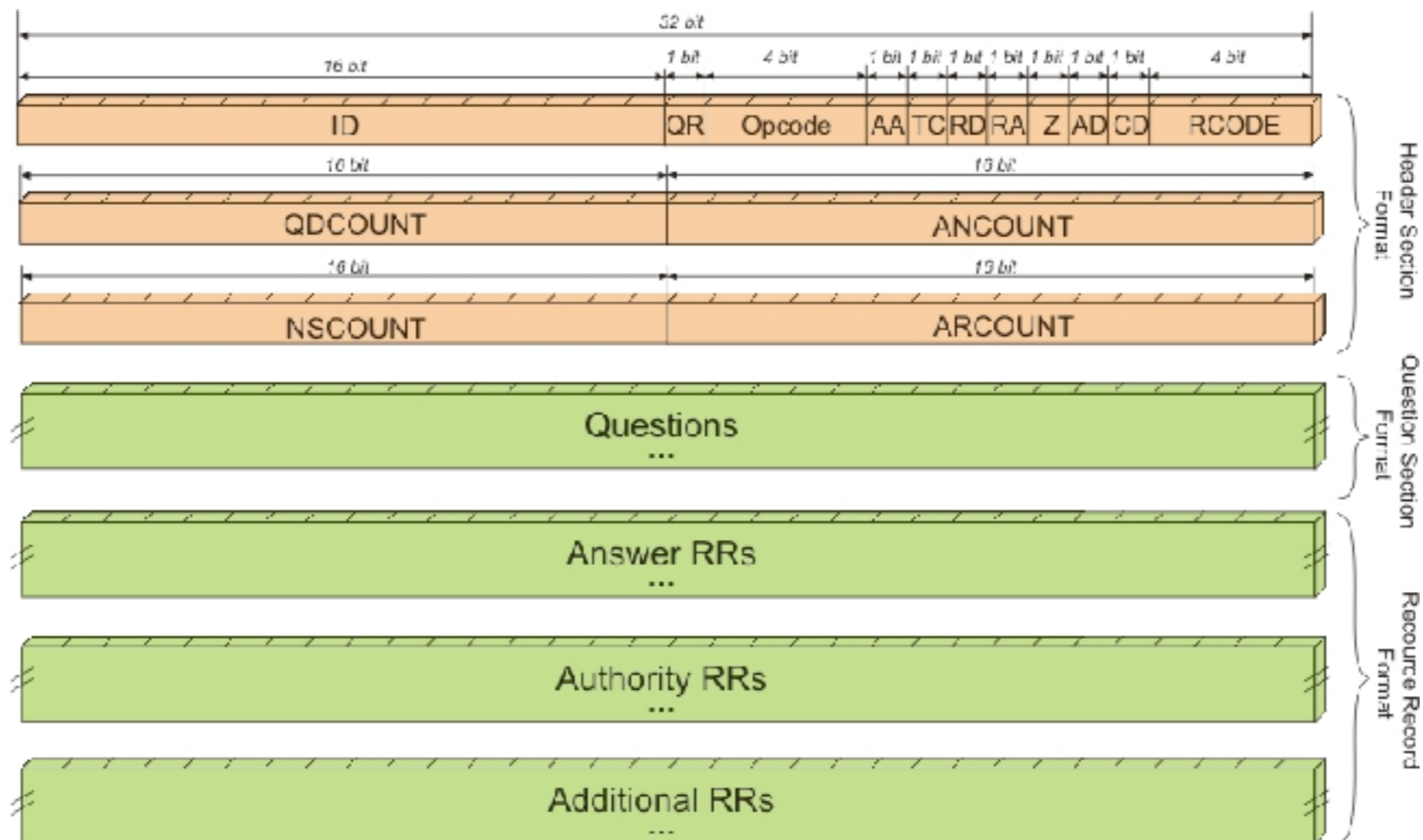# Deriving Correct Decoders



rewrite
MyDecoder!

```
SimpleDecoder (b : ByteString) :=
  (n, b) ← decodeNat(b);
  (s, b) ← decodeString(b); ???
```

⊔⊔

```
SimpleDecoder (b : ByteString) :=
  (n, b) ← decodeNat(b);
  (s, b) ← decodeString(b);
  (n', b) ← decodeNat(b);
  if (n' < 32) then
      (rs, b) ← decodeList(b, n);
      return ⟨ID :: s, readings ::rs⟩
else Error
```

# Parsing DNS Packets

- Synthesized decoder for DNS Packets (RFC 1035)
  - Specification ≤ 110 LOC
  - Valid: Compressed + Uncompressed packets
  - Data-dependent behavior used to parse response sections
  - Variable-type resource records

Narcissus in Action

# Evaluation

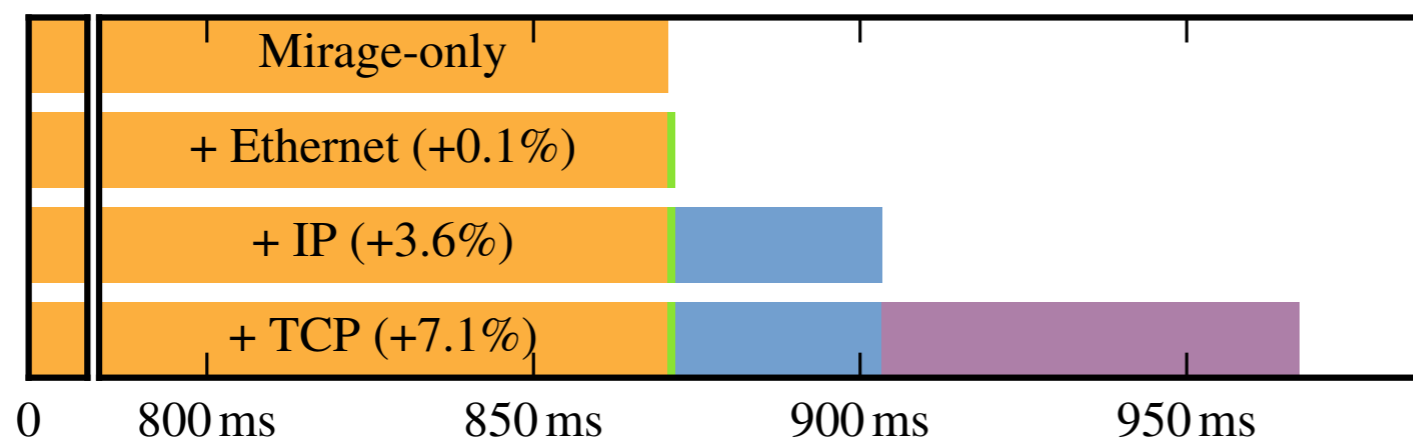| Protocol | LoC | Interesting Features |
|----------|-----|---------------------|
| Ethernet | 150 | Multiple format versions |
| ARP | 41 | |
| IP | 141 | IP Checksum; underspecified fields |
| UDP | 115 | IP Checksum with pseudoheader |
| TCP | 181 | IP Checksum with pseudoheader; under-specified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

## Derived Decoders

| Format | LoC | LoP | Higher-order |
|--------|-----|-----|--------------|
| Sequencing (**ThenC**) | 7 | 164 | Y |
| Termination (**DoneC**) | 1 | 28 | Y |
| Conditionals (IfC) | 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

# Evaluation

- MirageOS is a library operating system for secure, high-performance network applications written in OCaml
- Replaced network stack of MirageOS with extracted OCaml implementations of synthesized decoders.
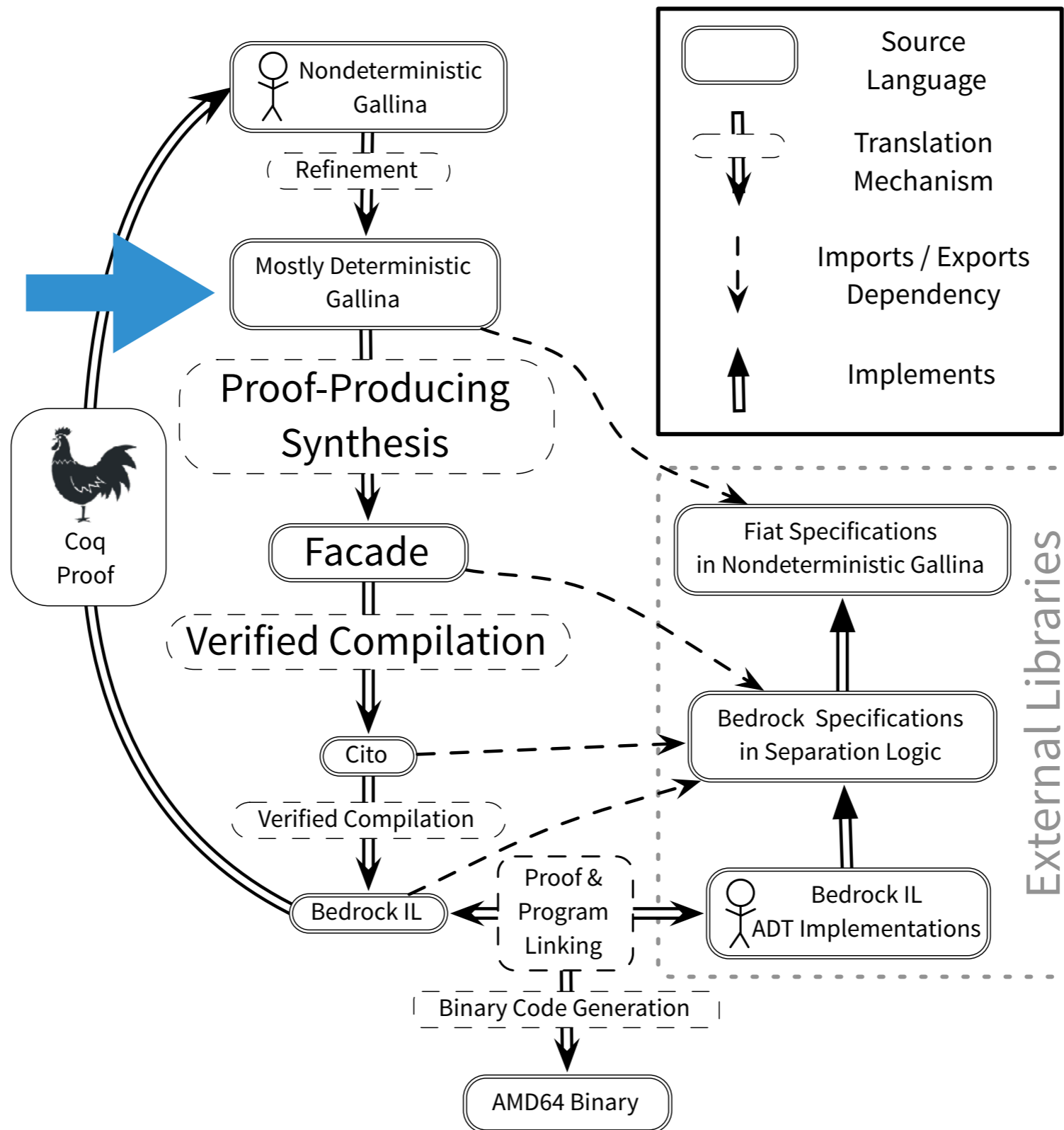- Found one problem in the test suite.



average page load time

performant?

# Synthesizing Performant Code

Fiat_Spec
Binary Format

Fiat_Impl
Mostly Deterministic
Functional Implementation

$\{(\text{Spec}, \text{Fiat}_{\text{Impl}}) \mid$
$\text{Fiat}_{\text{Impl}} \precsim \text{Fiat}_{\text{Spec}} \wedge$
$\{\text{Impl} \mid \text{Impl} \precsim \text{Spec} \wedge \text{Deterministic Impl}\}$
$\rightarrow \{\text{Fiat}_{\text{Impl}}' \mid \text{Fiat}_{\text{Impl}}' \precsim \text{Fiat}_{\text{Spec}}$
$\wedge \text{Deterministic Fiat}_{\text{Impl}}'\}$

Facade_Impl
Imperative Implementation

$\Gamma \approx \overline{\text{Spec}} \wedge$
$\Gamma \vdash \begin{array}{l} \text{meth}_i \mapsto p_i \mid \\ \forall \overline{v}. \ [\overline{x \leftarrow v}] \underset{\varnothing}{\overset{p_i}{\rightsquigarrow}} [\text{ret} \leftarrow \text{Fiat}_{\text{Impl}}.\text{meth}_i(v)] \end{array}$

Bedrock_Impl
Assembly Implementation

$\Gamma \vdash \begin{array}{l} \{\text{meth}_i \mapsto o_i \mid \\ \forall \Sigma. \ \{(\Sigma, p_i)\downarrow \wedge \ \text{state}(\Sigma)\} \quad o_i \\ \{\exists \Sigma'. \ (\Sigma, p_i) \Downarrow \Sigma' \wedge \ \text{state}(\Sigma')\} \end{array}$

Verified Assembly
Implementation of $\Gamma$

# The Future?

```
Definition DnsSchema :=
  Schema [ relation RECORDS has schema
              <NAME :: name,
               TTL :: nat,
               CLASS :: RRecordClass,
```
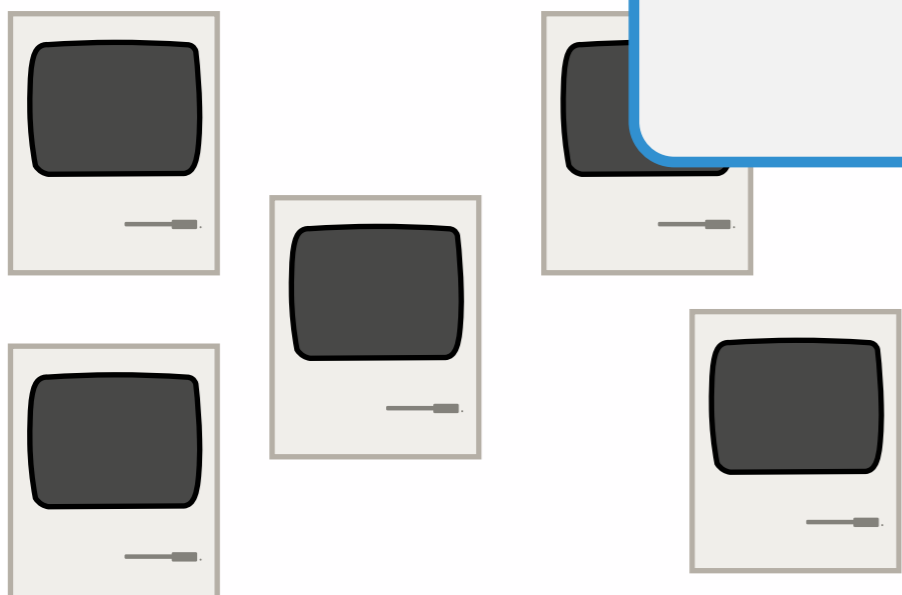
On-Demand L

On-Demand L

www.fa

31.13.69.2

```
ADT RRecordD
RepType := R
def DBInit()

def AddRecor
    Insert rr

def FindReco
    SortedBy
```

```
Definition PacketParserFormat :=
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|              ANCOUNT           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|              NSCOUNT           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|              ARCOUNT           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|
|
/              QNAME              /
/                                 /
/
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|              QTYPE             |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|              QCLASS            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+ …
```

```
ADT PacketParserSpec {
  def ParsePacket (request : ByteString) :=
      PacketParserFormat⁻¹ request,
  def EncodePacket (packet : Packet) :=
      PacketParserFormat packet,
    … }
```

ge();

request
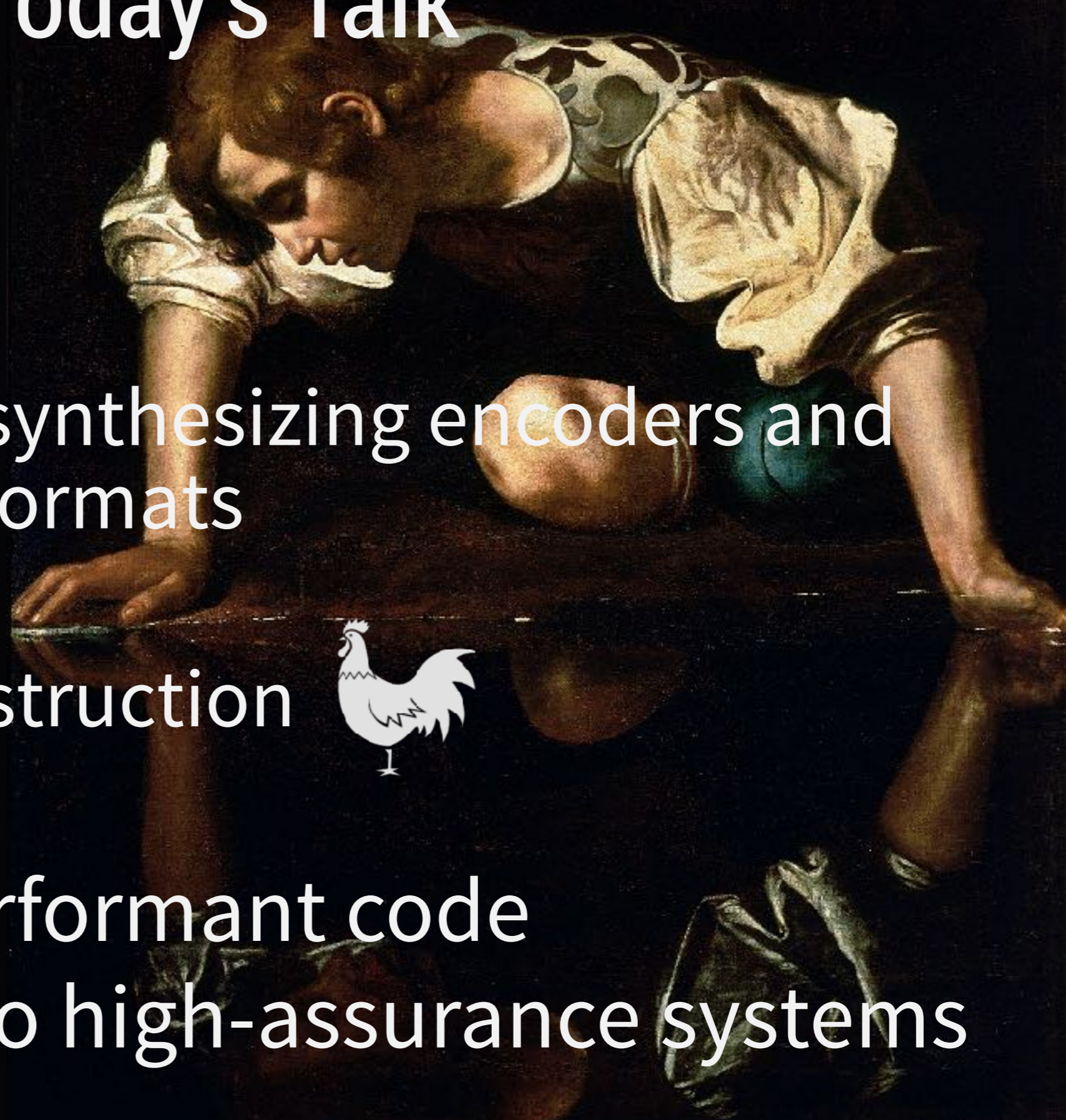
qname());

**fiat**

## Implemented in
## Coq Proof Assistant

- Rich higher-order logic
  for specifying program behavior

- Powerful tactic language for automating search
  for exploring implementation space

- Small trusted code base
  for certifying implementation meets specification

# Today's Talk

- **Narcissus:**
  - Framework for synthesizing encoders and decoders from formats
  - User extensible
  - Correct-by-Construction
- Generating performant code
- Integration into high-assurance systems