

# Smart Unpacking Research: Using Mathematics To Unpack More



Benjamin Long

**NIST** United States Department of Commerce  
National Institute of Standards and Technology

# Disclaimer

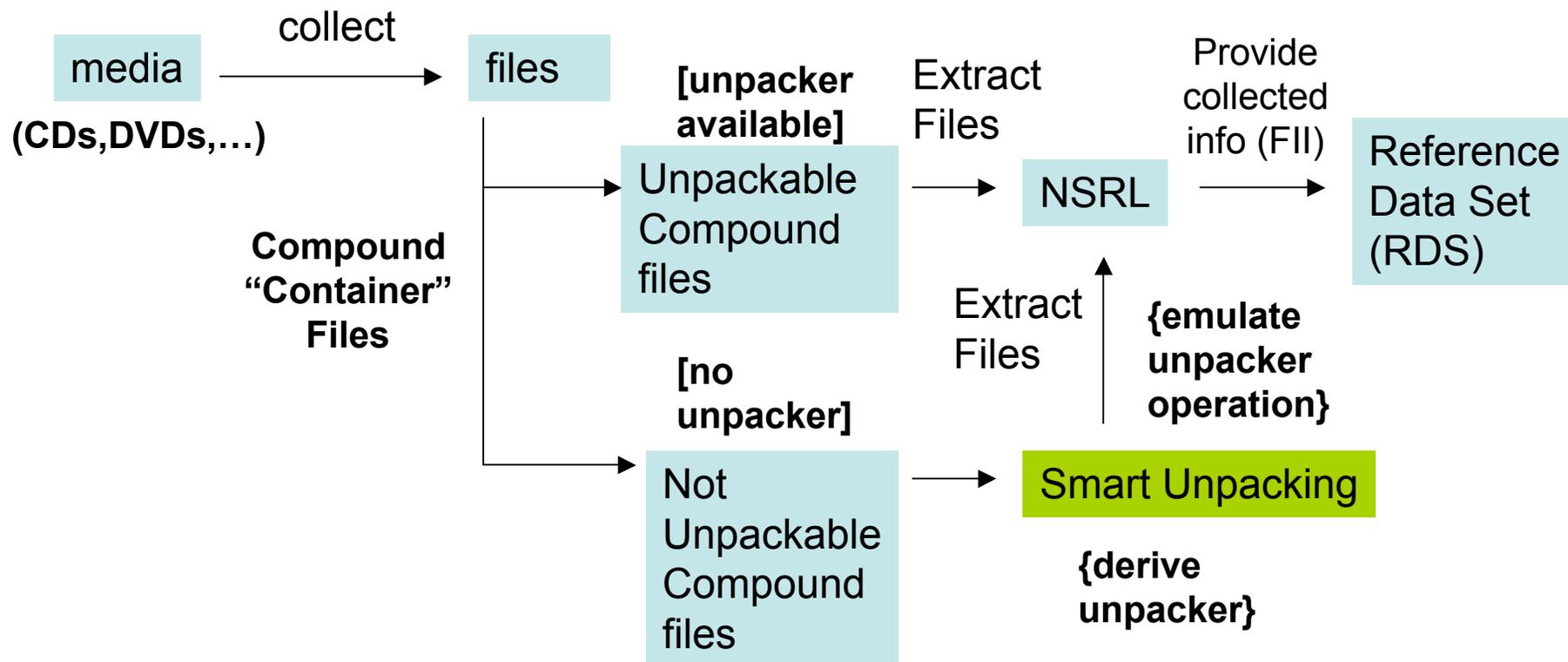
## **Disclaimer**

Trade names and company products are mentioned in the text or identified. In no case does such identification imply recommendation or endorsement by the National Institute of Standards and Technology, nor does it imply that the products are necessarily the best available for the purpose.

## **Statement of Disclosure**

This research was funded by the National Institute of Standards and Technology Office of Law Enforcement Standards, the Department of Justice National Institute of Justice, the Federal Bureau of Investigation and the National Archives and Records Administration.

# Context – Smart Unpacking and NSRL Operations



# Smart Unpacking

- NSRL collects files by unpacking
- New file types appear faster than new unpackers
- NSRL goal: unpack and provide FII for as many files as possible
- SU – attempts to facilitate this goal by
  - Unpacking files when have no unpacker
  - Forming measurements (completeness, accuracy) of unpacking ops

# Nature of Unpacking

- Unpacking: locate, identify, and extract files completely and accurately from compound structures
- SU studies and formalizes
  - How to use mathematical methods to find structure in content and nonstandard packaging
  - How to use this to detect, extract, and measure files

# SU Research – Early Stages

- To date has –
  - Established basis for theoretical and practical implementation
  - Not yet formalized unpacking systems
  - Implemented specific prototypes and experiments
- Experiments have shown capability to
  - Do basic inline extractions
  - Do similarity and grammatical pattern correlation across many file types (binary and non-binary)



# Technical and Scientific Basis and Methods

Methods used to Discern Structure and Metadata

-Modeling

-Uniqueness features in – syntax, content, abstract pattern representations

## Pattern Theory

$G = \{g_0, g_1, \dots, g_i, g_i^0\} = \text{generator\_space}$

$S : G \longleftrightarrow G, s \in S = \text{similarity\_group}$

$G = \bigcup_{\alpha \in A} G^\alpha = \text{partition\_of\_generators}$

$b_1, b_2, \dots, b_\omega = \text{bond\_values}$

$B_s(g) = \{b_j; j = 1, 2, 3, \dots, w(g)\} = \text{bond\_structures}$

$B_v(g) = \{\beta_j = 1, 2, \dots, w(g)\}$

...

## Formal Language Theory

### Context-free grammars

$G = (V, \Sigma, R, S) = \text{grammar}$

$V = \text{nonterminals}$

$\Sigma = \text{terminals}$

$S = \text{start\_rule}$

$R = \text{rules}$

...

## Mathematical Modeling based on:

- Information Theory
- Statistics
- Probability
- Universal Algebra
- Topology
- Graph Theory, and more ...

## Parser Theory

- Variable lookahead mechanisms
- Multi-channel token processors
- Augmentation with syntactic and semantic predicates for context-sensitivity

## Measurements and Measurability

- Derived from mathematical models

# Spectrum of Unpacking Scenarios

- Simplest – very little knowledge req'd, inline metadata and content; just need an unpacker
- More complex – more complicated structure, nesting, intermingled file, content, and encoding types, alignment considerations (MSI experiment)
- Harder – req's knowledge of metadata and structure relationships; mathematical + experimental models
- Even Harder – at least 1 transform – compressed, encrypted, etc – but still can get metadata and op properties
- Hardest – 2 or more transforms, little to no metadata – will often still get some metadata and specific properties

# Unpacking Operation Measurement

- Measuring (at least) – accuracy and completeness of unpacking operation
  - Per file – from a given compound file
  - Across all files/content from a given compound file
- Beginning – focus on simple, coarse-grained ratios – total found/total extracted (files, content)
- Ongoing Objective –
  - develop content and format-specific measures – entropy-, probabilistic-, algebraic-based properties of content and structure
  - Simplify into classes of overall properties per file type

# Work From Here

- Integrate prototype tools and experiments into automated unpacking framework
- Research and integrate
  - methods for harder cases
  - Relevant measurements
- Track ability to increase overall coverage of the NSRL

# Contact

**Benjamin Long**  
**[www.nsrl.nist.gov](http://www.nsrl.nist.gov)**  
**[nsrl@nist.gov](mailto:nsrl@nist.gov)**

**Barbara Guttman**  
**Software Diagnostics & Conformance**  
**Testing Division**  
**[barbara.guttman@nist.gov](mailto:barbara.guttman@nist.gov)**

**Sue Ballou, Office of Law Enforcement**  
**Standards**  
**Rep. For State/Local Law Enforcement**  
**[susan.ballou@nist.gov](mailto:susan.ballou@nist.gov)**