



# **Grammar Inference via Dynamic Taint Tracing**

Pavel Zakopaylo, ANU



# Motivation

- Automated security analysis of software applications (i.e. finding bugs)
- Want to describe the input data format;
- Application is given as a compiled binary – no source code or debugging symbols
- Analysis of data flow in memory is important to achieving these aims



# Prior Work

- TUPNI – Microsoft Research
- Taint Tracing as Data Flow Analysis
- Taints :: memory location  $\rightarrow$  { file offset }  
... at any point during execution
- TUPNI paper does not precisely describe how these taints are determined, and cites papers that are similarly vague



# My Taint Tracer

- Built using Intel's PIN
- Intercepts Linux syscalls on IA-32 and Intel-64 architectures
- Inserts update to taint database after each (relevant) instruction
- At attempt to replicate the undocumented features of TUPNI

# Implementation

- PIN's API provides some functionality to determine which memory locations / registers are read / written
- Mostly dealing with edge cases
- e.g. PUSH, POP affect RSP deterministically, so do not spread taint
- Toy parsers written as „unit tests“

# Performance

- Most (non-TUPNI) work on taint tracing cites „negligible“ performance overhead
- Most likely b/c they only track *if* a location is tainted (i.e. memory overhead  $\leq 12.5\%$ )
- My testing shows **orders of magnitude** more memory & CPU usage
- Performing an  $O(\log N)$  operation after every instruction is **slow**



# // TODO

- Mitigate performance issues, to the point where it's at least usable
- Make tool connect to an actual grammar inference system
- (Or to an architecture-independent format as originally suggested by Shane)
- More special cases
- Replicate TUPNI (!!)



# Thank You

- Questions?