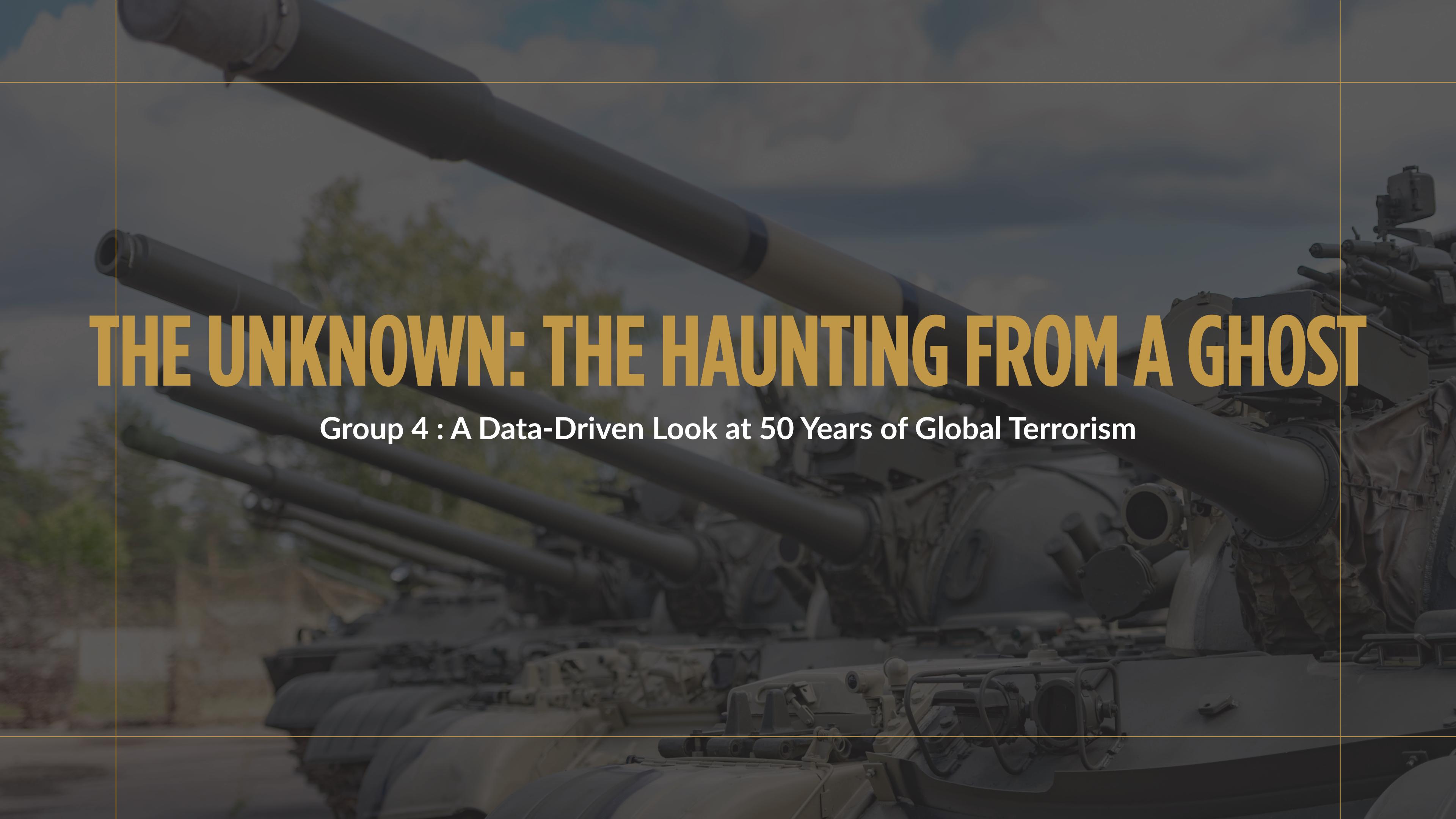


# **THE UNKNOWN: THE HAUNTING FROM A GHOST**



**Group 4 : A Data-Driven Look at 50 Years of Global Terrorism**

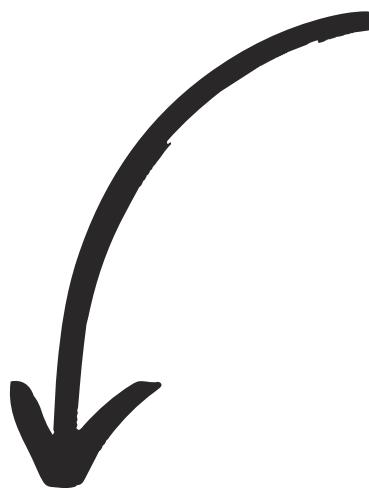


# OPERATION GHOST HUNT: THE INVESTIGATION

Our mission is to transform **chaotic data** into **actionable intelligence** by answering two critical questions:

## 1. CLEANSE THE CRIME SCENE: Data Preparation

How many "ghosts"—missed signs or miscoded warnings—are hidden in the raw data, and how do we purge them?



## 2. PORTRAIT OF THE INVISIBLE ENEMY: Clustering

After uncovering these ghosts, what do they reveal about the true faces of the enemy we are fighting?





**FEAR IS LOUD. DATA IS QUIET.**

We chose to listen to the data

## OUR INVESTIGATION

### SOURCE

The Global  
Terrorism  
Database

### SCOPE

1970 - 2017

### FOCUS

Systemic truths,  
not singular  
events

01

02

03

04

05

06

07

# THE DATA CHALLENGE - THE "FOG OF WAR"

181,691 Raw Incidents    135 Features

## MISSING INTELLIGENCE

45% of features are >90%  
empty

## INCONSISTENT CODES

"Unknown" has many names:  
-9, -99, NaN, "Unknown", etc

## ILLOGICAL RECORDS

The Data Contradicts Itself:  
hours kidnapped < 0

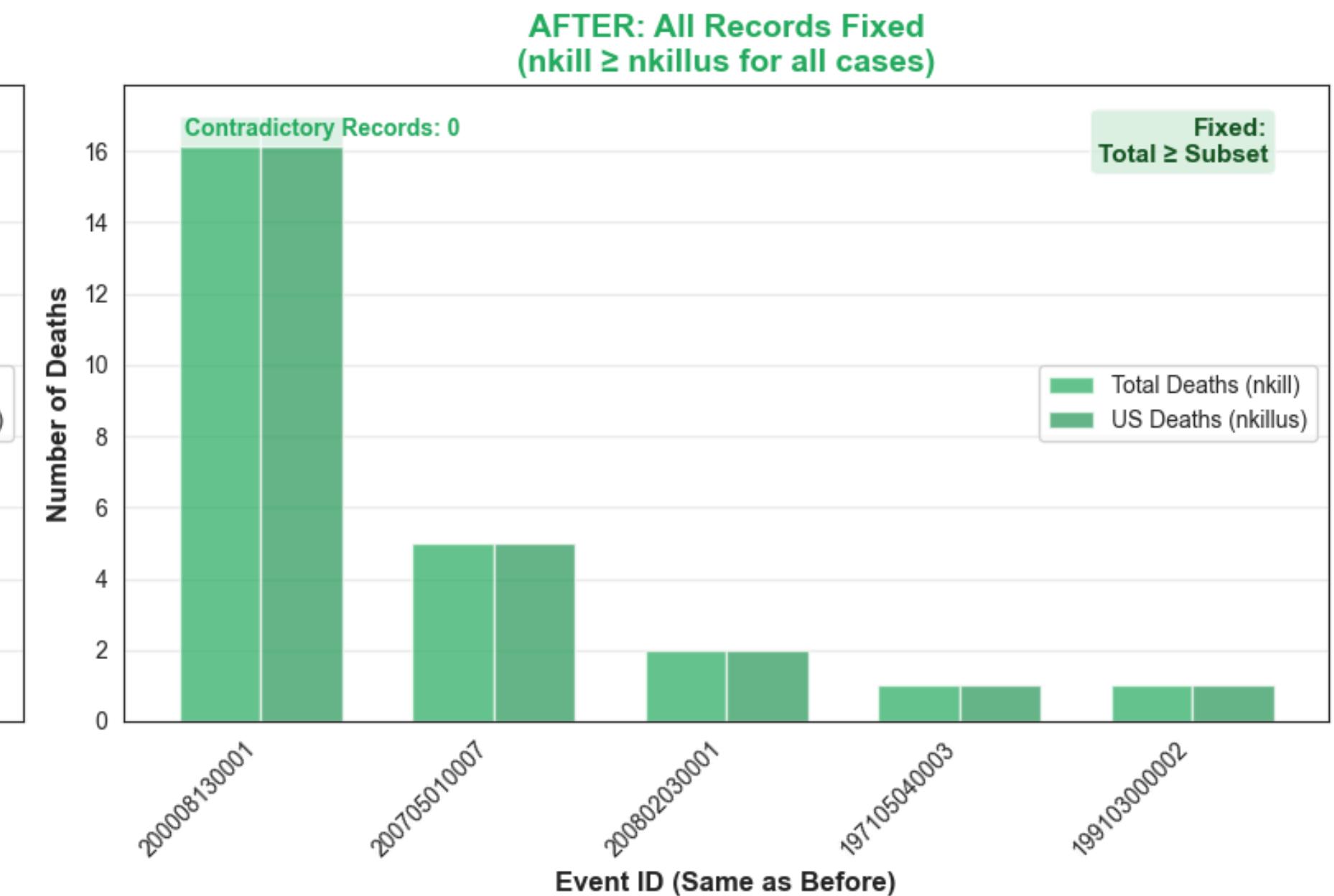
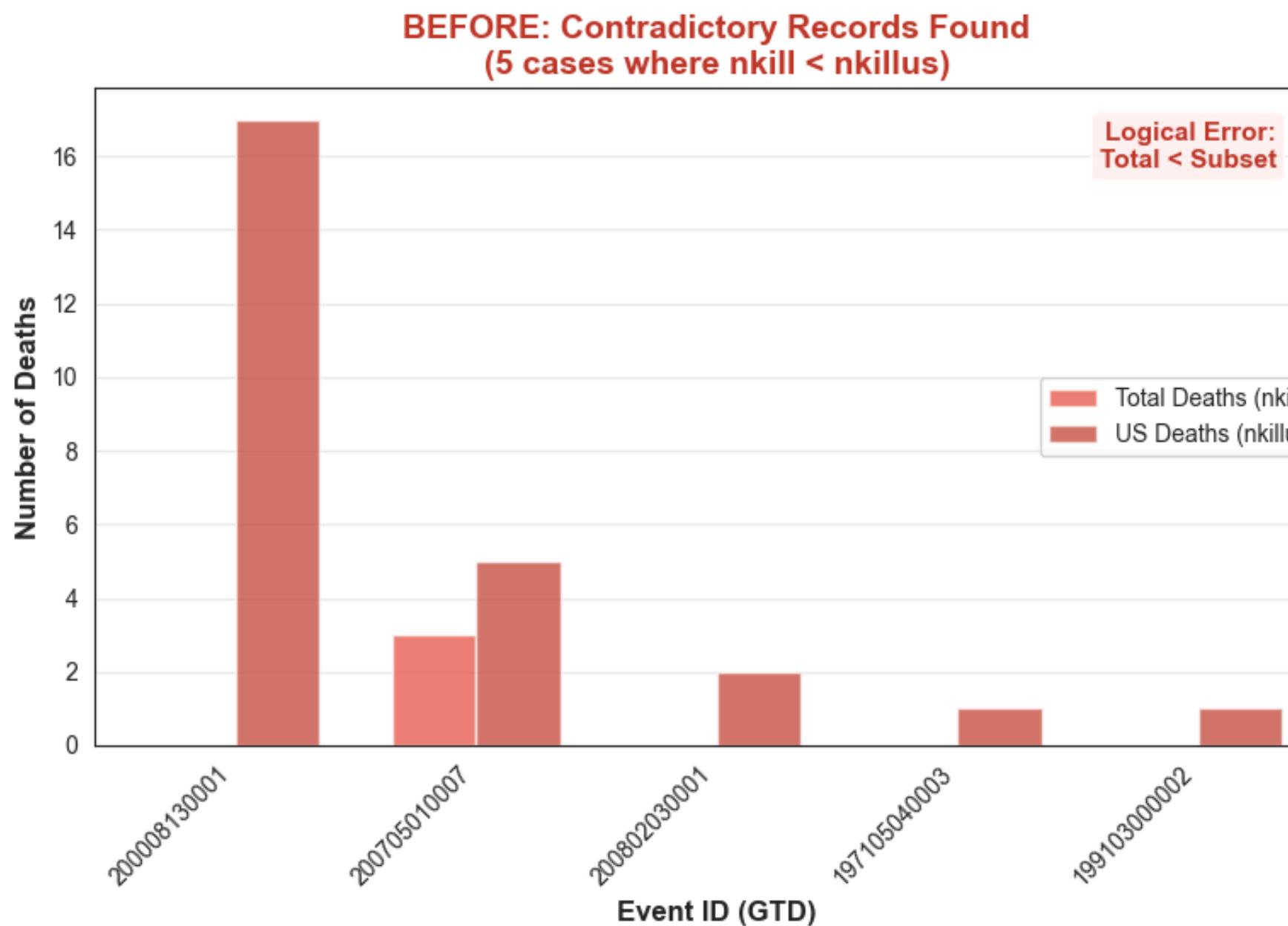
## INVALID INFORMATION

Impossible Dates: 0/5/1985,  
14/0/2004

# THE FIRST CLUE: CONTRADICTORY STATEMENTS

The Logical Error: Total Casualties cannot be less than U.S. Casualties ( $nkill < nkillus$ )

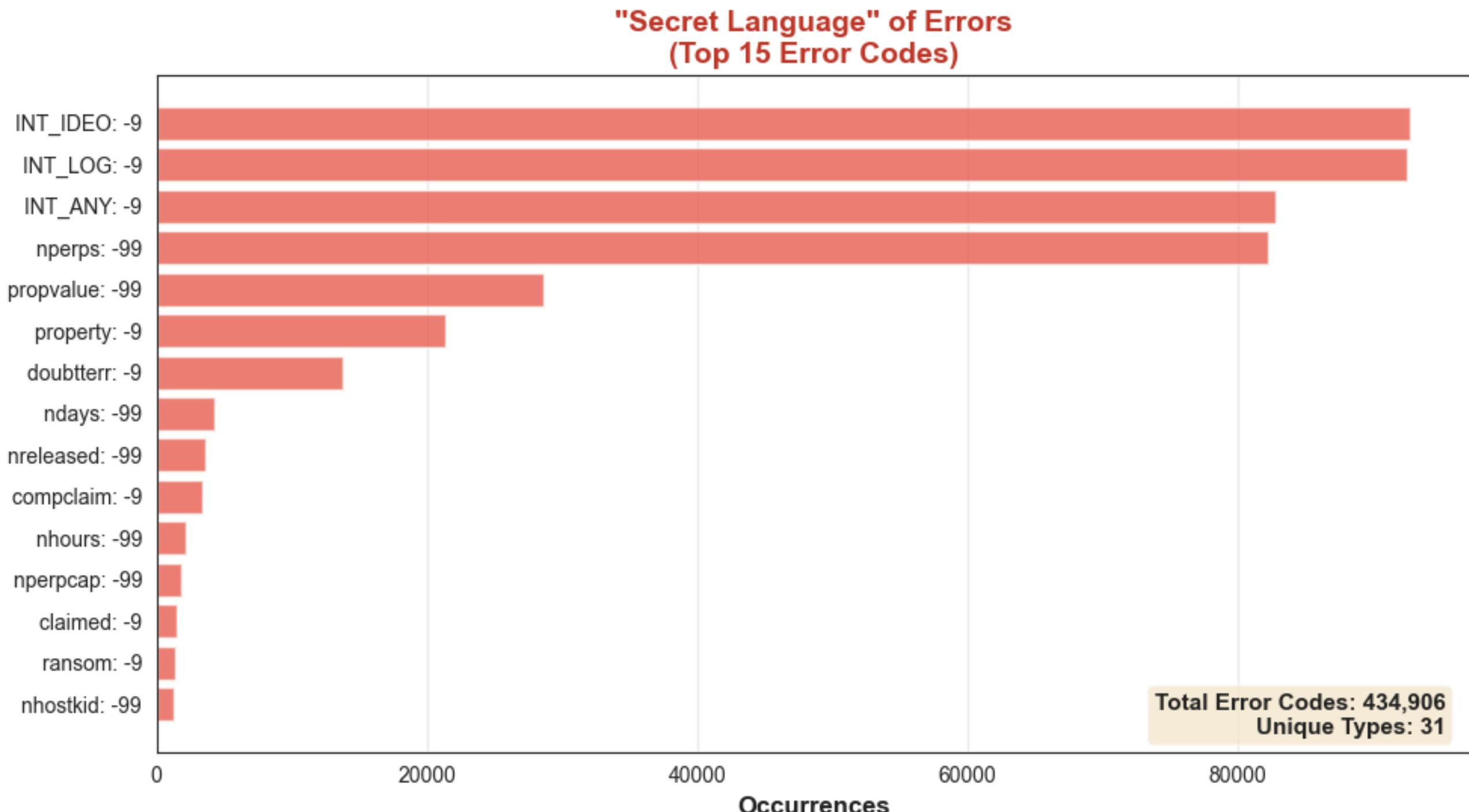
First Clue - "Contradictory Statements" Fixed



# THE SECRET LANGUAGE OF OMISSIONS

Standardizing **-9**, **-99**, **NaN**, and **Unknown** into a Unified Term.

## Second Clue - "The Secret Language of Missing Data"



Before

**31 DIFFERENT ERROR CODE TYPES**

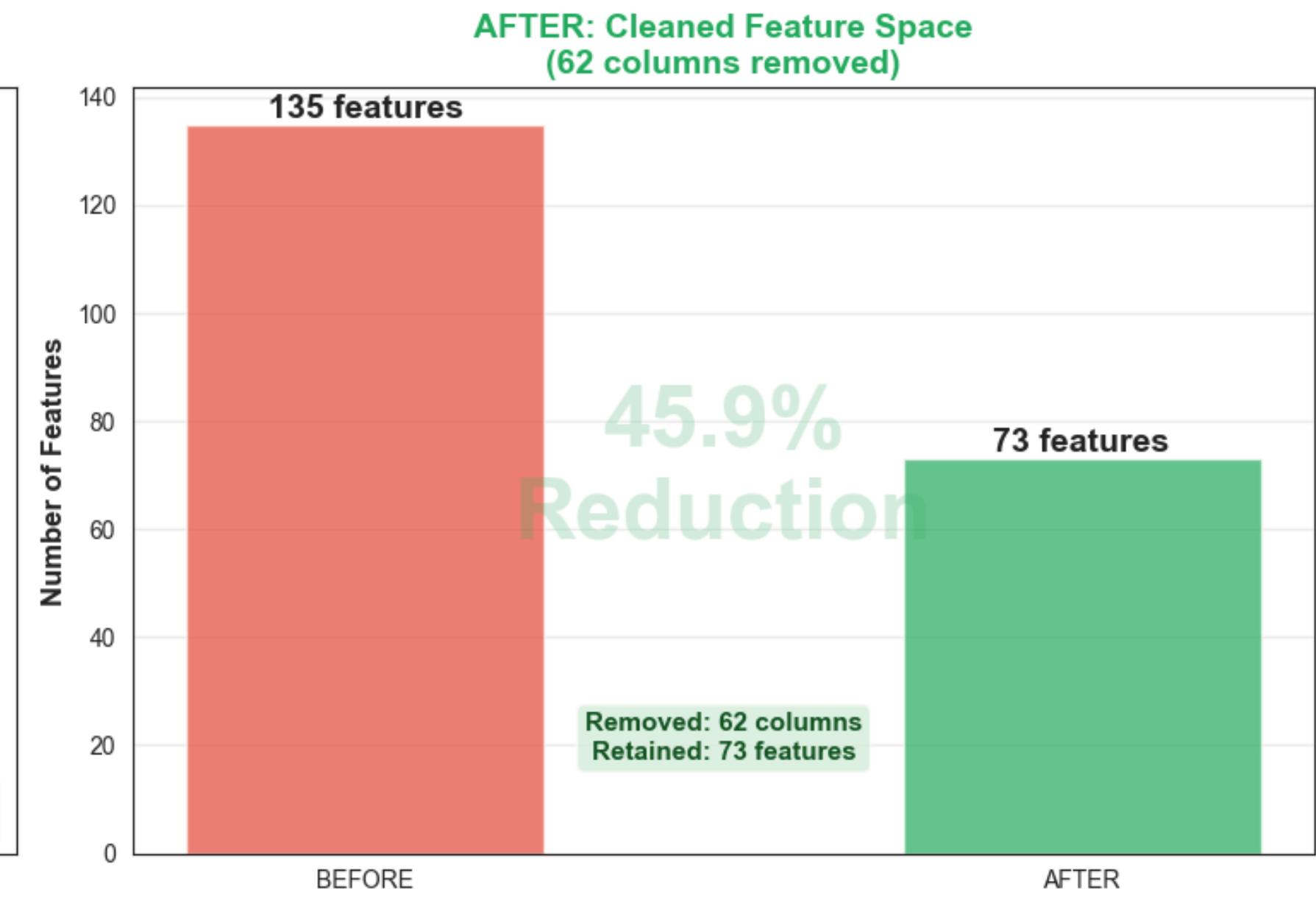
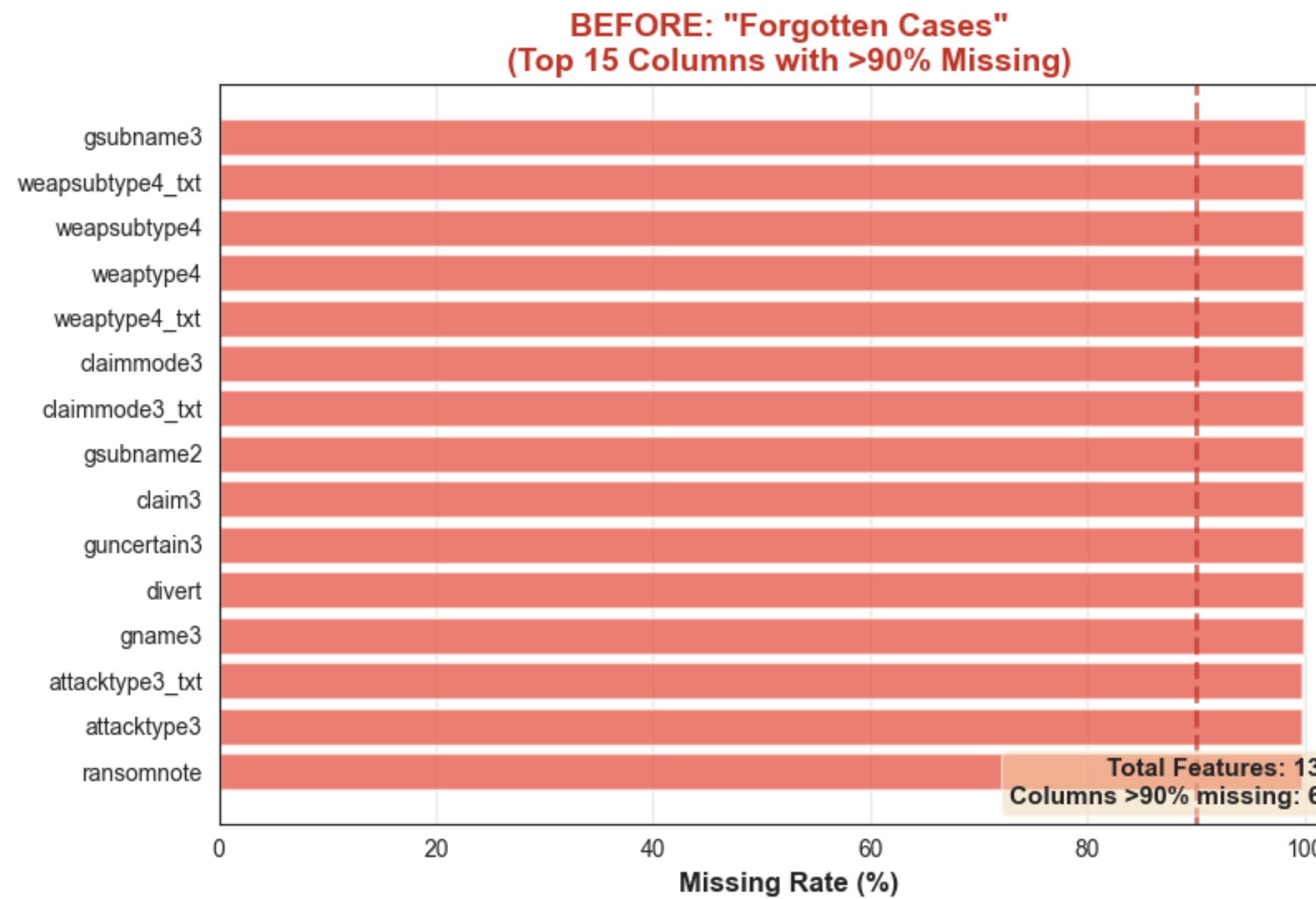
After

**1 UNIFIED (NAN)**

# THE THIRD CLUE: THE FORGOTTEN CASES

Feature Reduction: Eliminating Columns with >90% Missing Data.

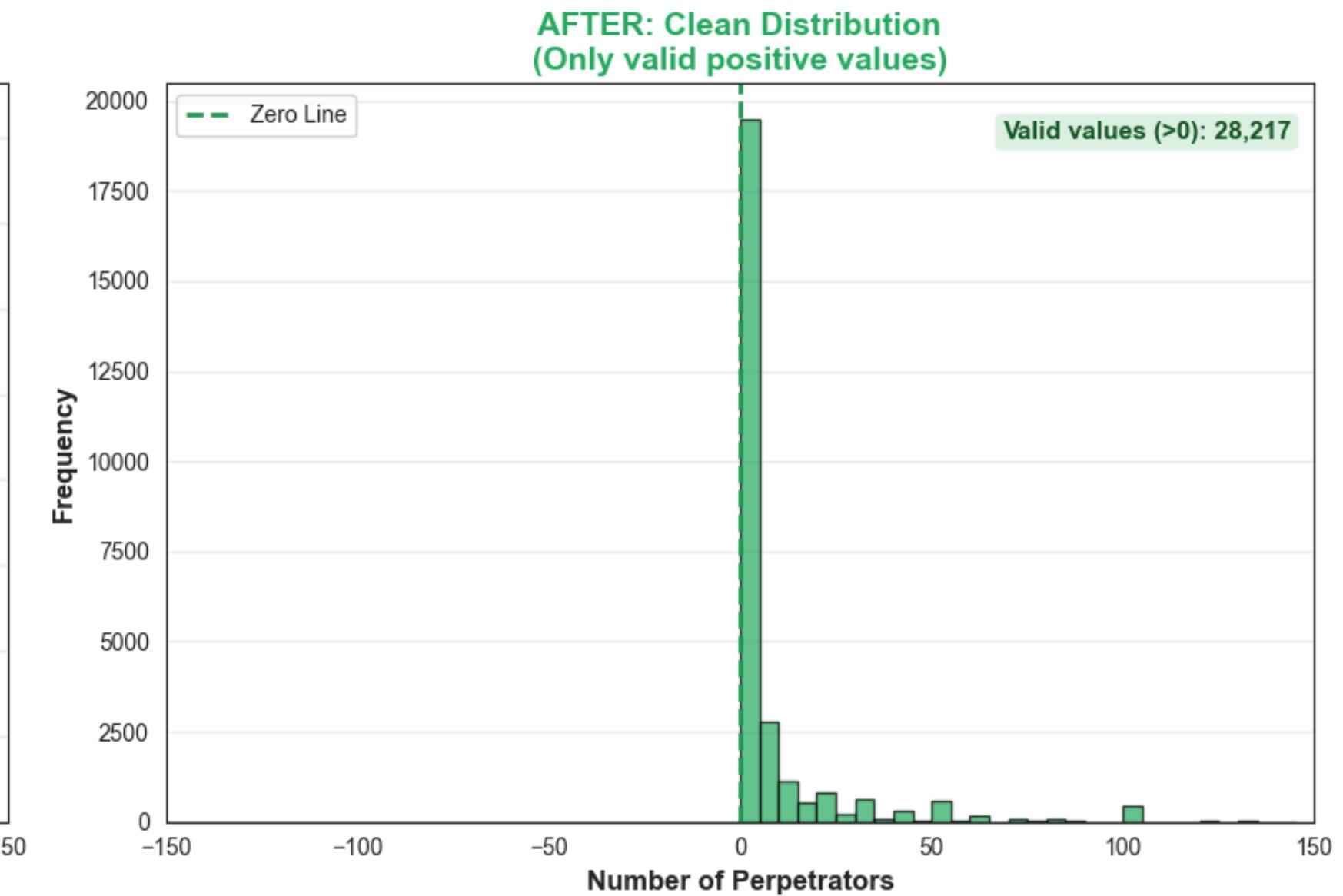
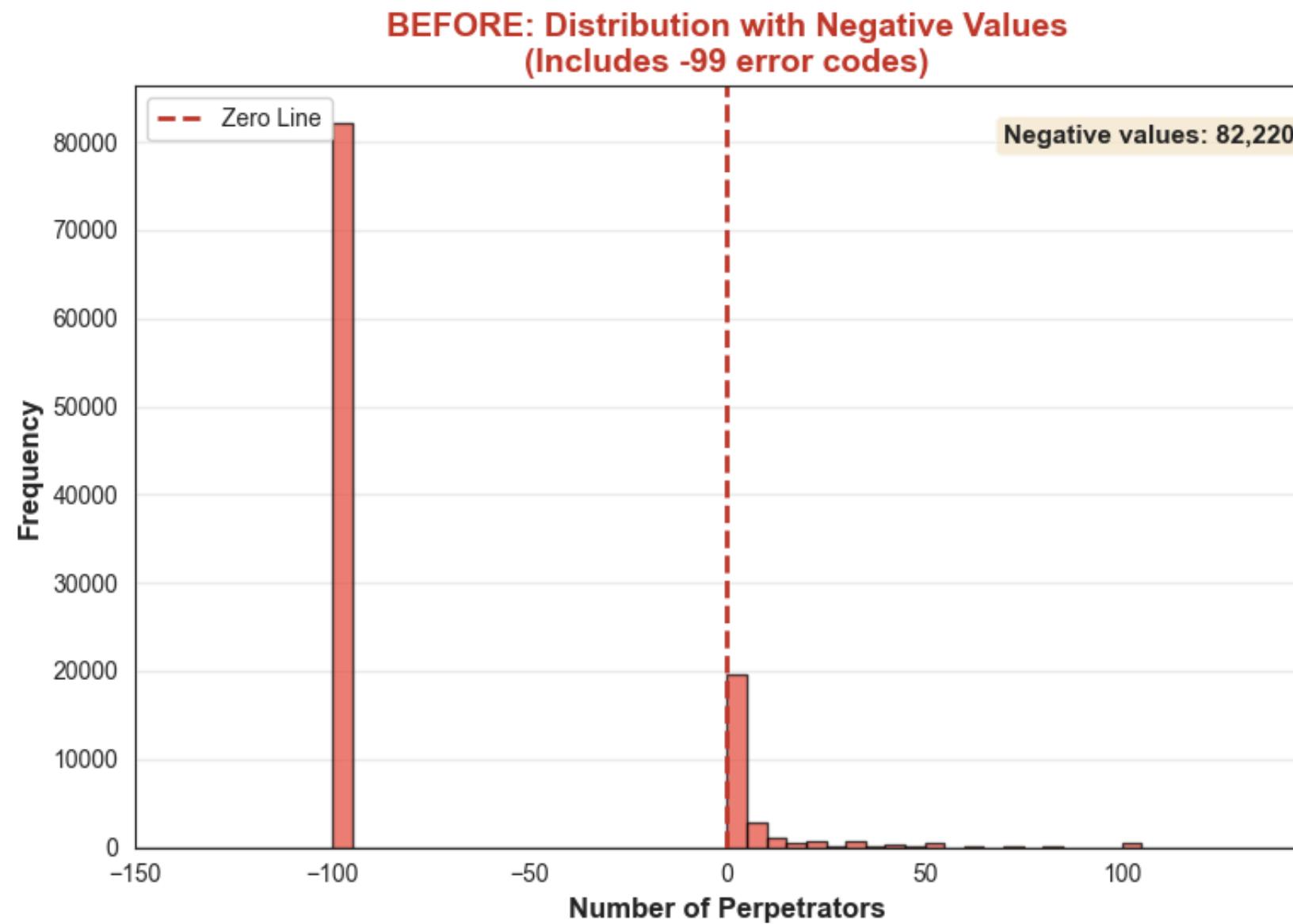
Third Clue - "Forgotten Cases" Removed



# CASE STUDY 1: EXORCISING GHOST VALUE FROM NPERPS

Handle Errors Outliers: Replace **non-positive** with **NaN**

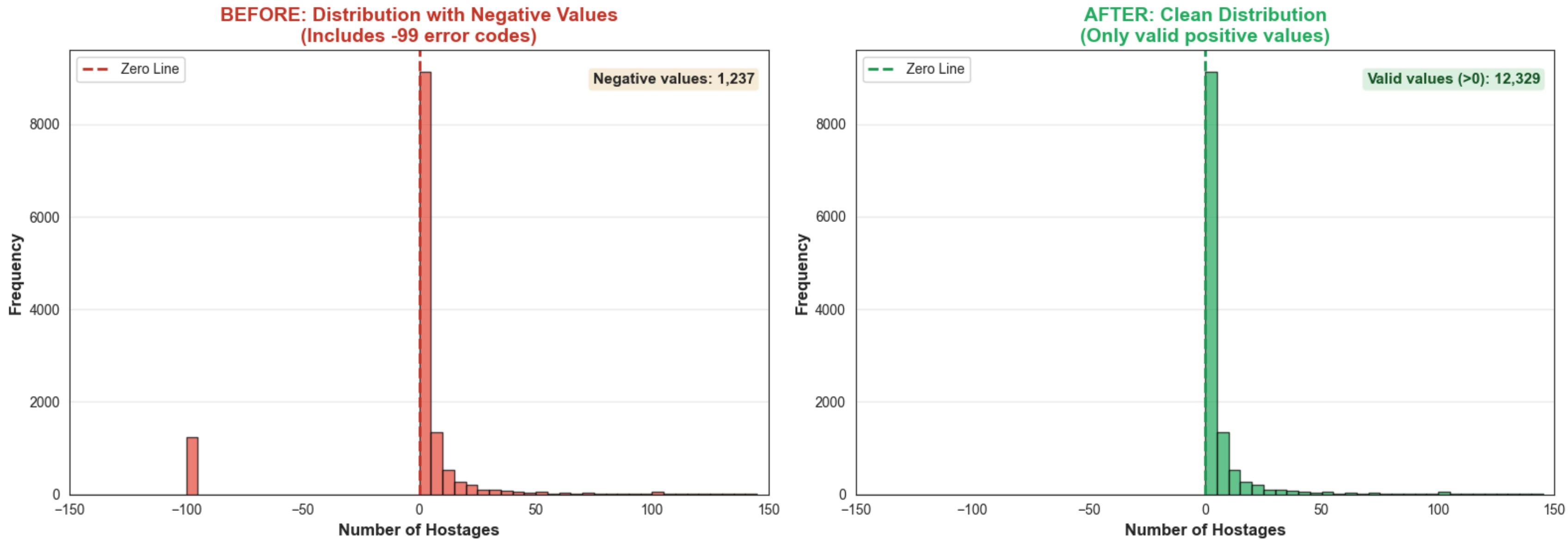
Case Study #1 - Cleaning nperps Column



# CASE STUDY 2: EXORCISING GHOST VALUE FROM NHOSTKID

Handle Errors Outliers: Replace **non-positive** with NaN

Case Study #2 - Cleaning nhostkid Column





# FROM CHAOS...

**135+ TOTAL FEATURES**

A sea of noise

**45% OF FEATURES  
>90% EMPTY**

Unreliable & incomplete

**5+ "MISSING" CODES**

Ambiguous, error-prone

**1000+ ILLLOGICAL RECORD**

Fundamentally broken

# TO CLARITY

**70 RELEVANT FEATURE**

Signal over noise

**0 MISSING VALUES (NAN)**

Complete for Analysis

**CONSISTENT "UNKNOWN" LABEL**

Clear & machine-readable

**0 LOGICAL CONTRADICTIONS**

Internally consistent & trustworthy



01

02

03

04

05

06

07

# THE GHOST HUNT: DRAW A PORTRAIT OF INVISIBLE ENEMY

**“UNKNOWN”**

THIS ISN'T AN ANSWER.  
IT'S AN **INTELLIGENCE FAILURE**

**OUR APPROACH**

INSTEAD OF **WHO THEY ARE**  
WE ASKED **HOW THEY ACT**  
TO FIND PATTERNS

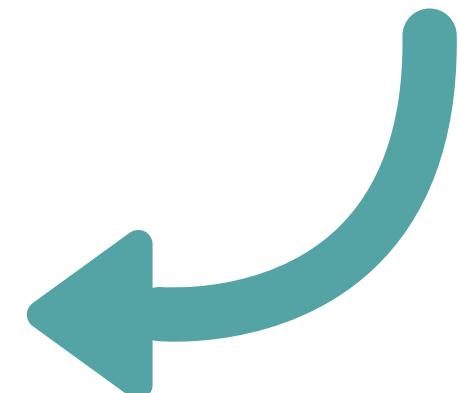


## WEAPON

MACHINE  
LEARNING



HUNT THE GHOST  
HINDING IN THE  
DATA



# THE FIRST ATTEMPT: CLUSTERING ON RAW DATA WHEN THE MACHINE TRIES TO READ A CORRUPTED FILE, IT SEES NONSENSE.

**Cluster Size Distribution**

153329



**Average Casualties per Cluster**



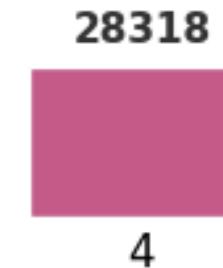
4 K:3.4 | W:4.6

1 K:0.2 | W:5.8

0 K:2.0 | W:2.4

2 K:0.0 | W:0.0

Avg Killed  
Avg Wounded



A **BROKEN** RESULT: **99%** OF THE DATA WAS FORCED INTO A SINGLE GROUP WHILE OTHERS ARE **EMPTY**

**ILLOGICAL** PROFILES: CLUSTER 3 SHOWS AN IMPOSSIBLE AVERAGE OF OVER **7,000 CASUALTIES PER ATTACK**

# **THE BREAKTHROUGH: FIVE GHOSTS IN THE MACHINE**

## **WITH CLEAN DATA, THE MACHINE COULD FINALLY SEE. THE "UNKNOWN" IS NOT ONE ENTITY - IT IS FIVE.**

### **PROFILE 1: THE AGITATOR**

The Low-Level Noise

### **PROFILE 2: THE GUERRILLA**

The Hybrid Fighter

### **PROFILE 3: THE SPECIALIST**

The Kidnapping Professional

### **PROFILE 4: THE MASS-CASUALTY BOMBER**

The Hyper-Lethal Threat

### **PROFILE 5: THE ASSASSIN**

The Precision Striker

# THE LANDSCAPE: MAPPING THE VOLUME OF THE GHOSTS

NEARLY 50% OF ALL 'UNKNOWN' INCIDENTS ARE LOW-IMPACT 'NOISE'.

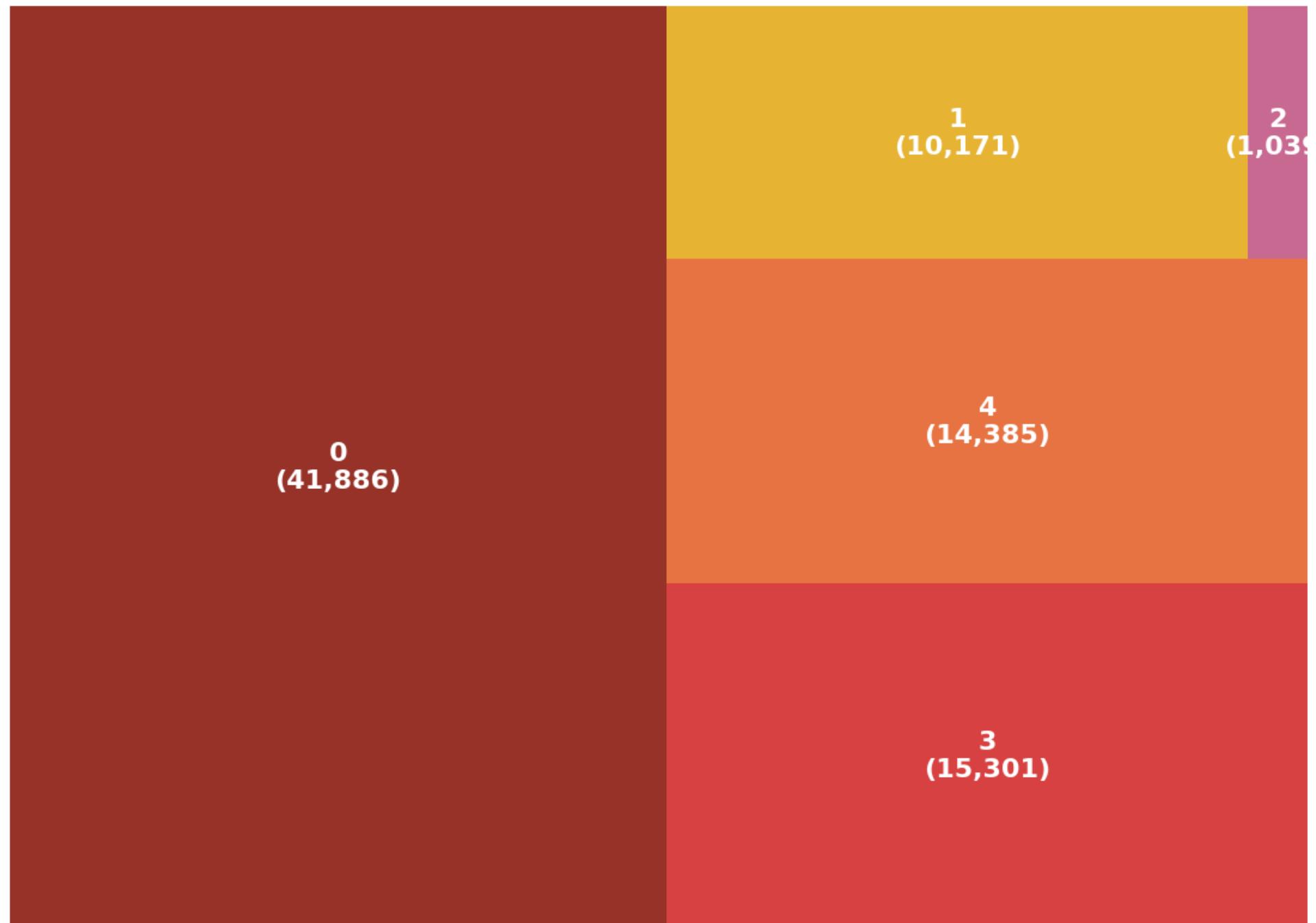
## CLUSTER 0 - THE AGITATOR

The dominant force by volume, creating the most 'noise'.

## CLUSTERS 3 - THE MASS-CASUALTY BOMBER

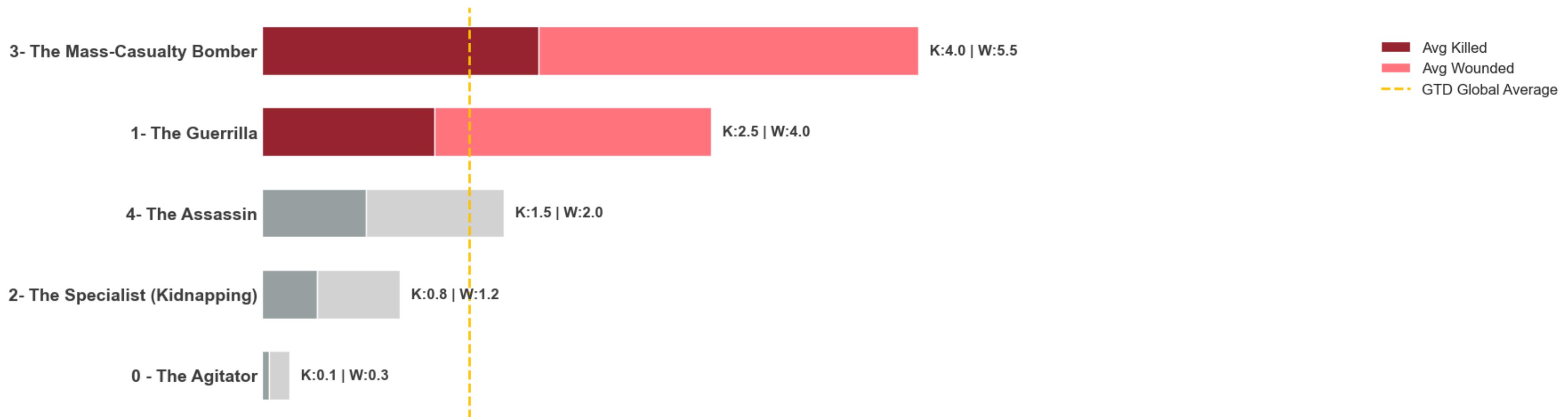
The second most frequent, representing a significant volume of high-risk attacks.

Volume: Half the attacks are from the 'Amateur' group



# THE MASS-CASUALTY BOMBER & THE GUERRILLA: OUR DEADLIEST INVISIBLE THREATS

## THE MASS-CASUALTY BOMBER & THE GUERRILLA: OUR DEADLIEST INVISIBLE THREATS.



**50% OF THE DATA (CLUSTER 0 - THE AGITATOR) IS  
LOW-IMPACT 'NOISE'**

they are resource-draining 'noise', not  
strategic threats.

**CLUSTERS 3 AND 1 REPRESENT THE ACTUAL  
KILLING MACHINES**

all counter-terrorism efforts must be re-  
focused on these two profiles.

# EACH GHOST HAS A UNIQUE TACTICAL FINGERPRINT

## THEIR METHODS ARE NOT RANDOM; THEY ARE SPECIALIZED. THIS IS THEIR OPERATIONAL DNA.

### CLUSTER 2 - THE SPECIALIST

Over 90% of this group's activity is a single tactic: Hostage Taking (Kidnapping).

They are professionals, not generalists.

### CLUSTER 4 - THE BRUTE FORCE

The Mass-Casualty Bomber relies almost exclusively on one method: Bombing/Explosion.

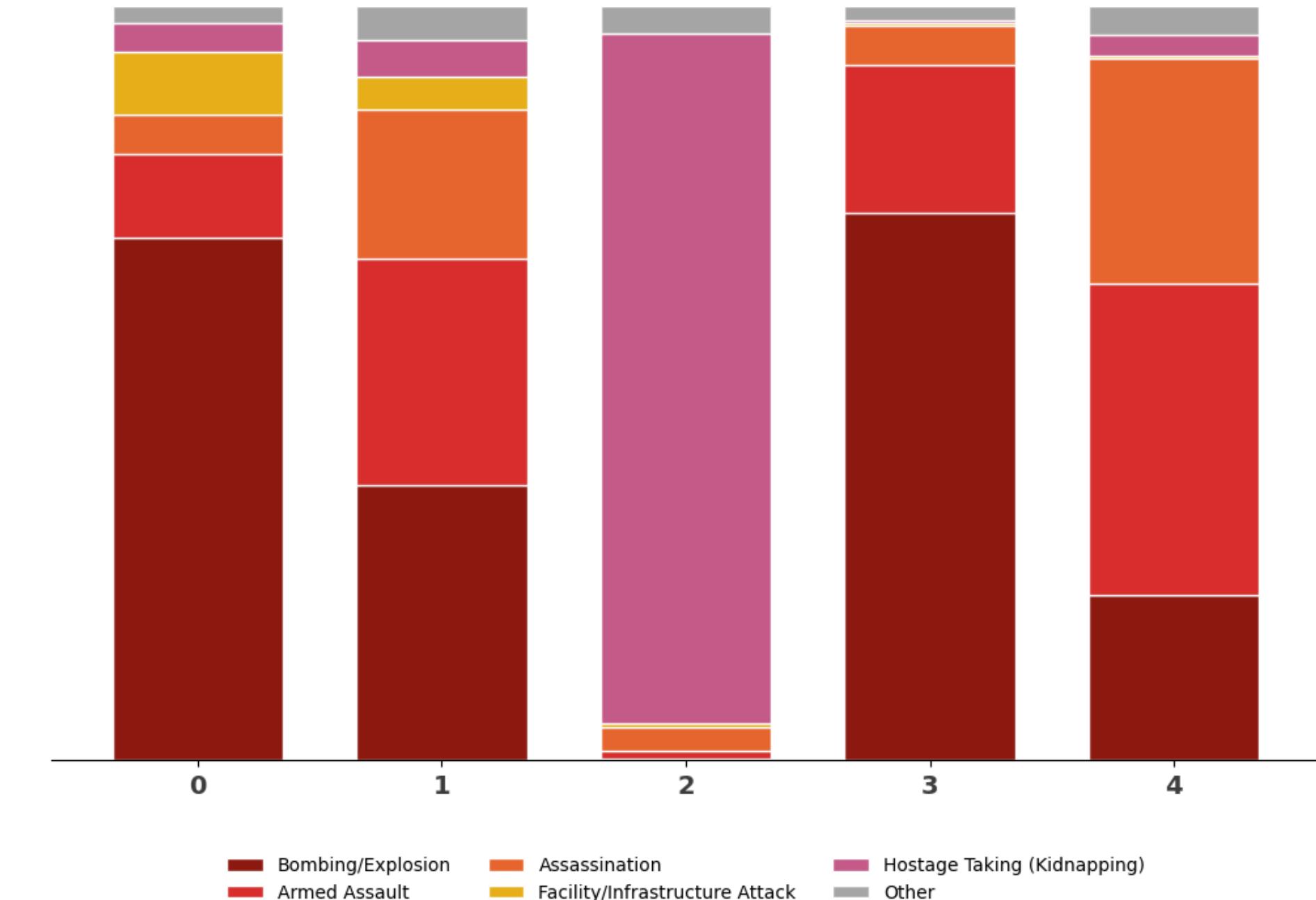
This is a high-impact, low-complexity strategy.

### CLUSTER 1 - THE HYBRID FIGHTER

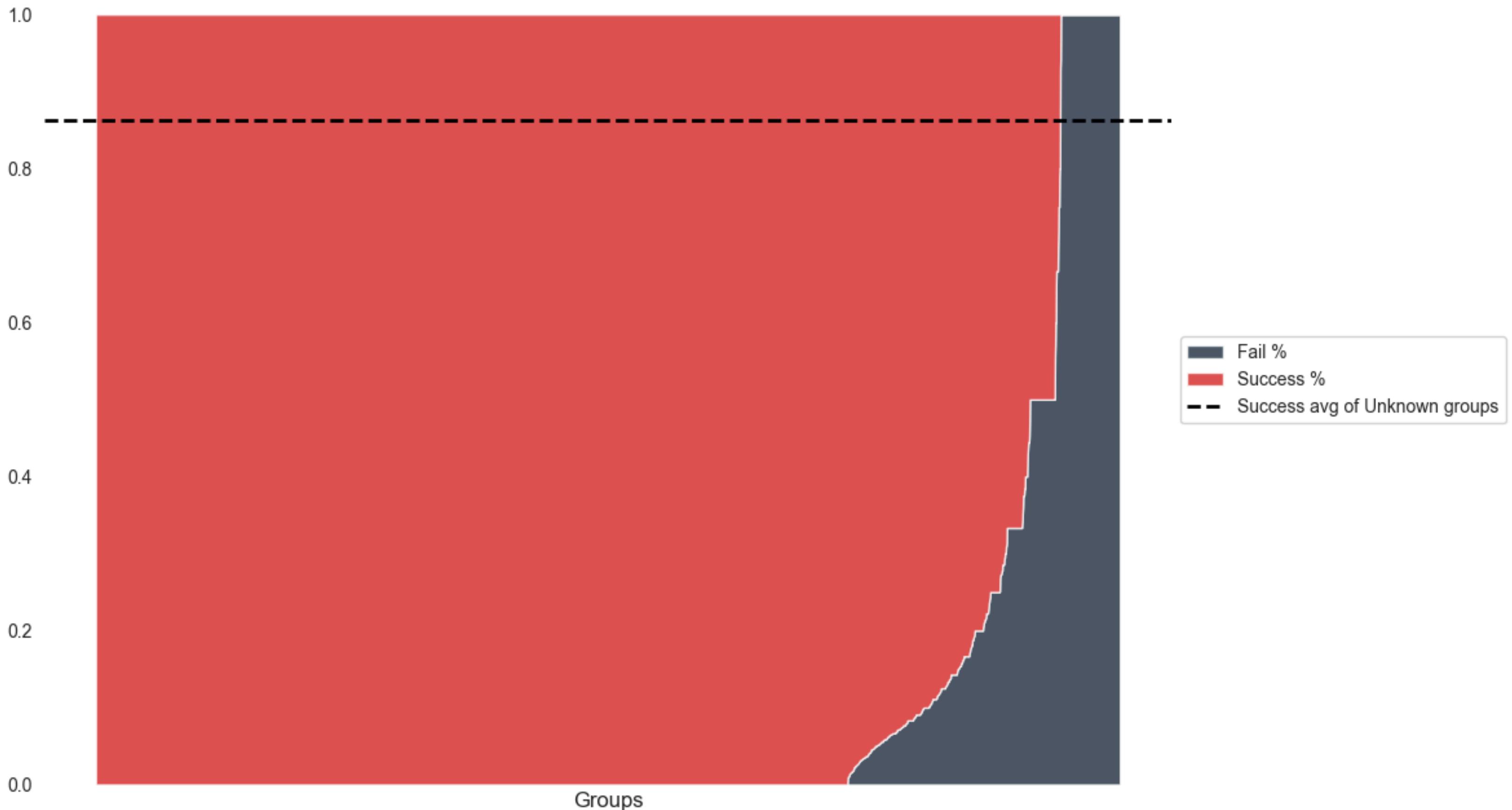
The Guerrilla shows the most tactical diversity, balancing Bombing, Armed Assault, and Assassination.

This suggests a more organized, flexible adversary.

Tactics: Most use Armed Assault and Bombing



# DO NOT MISTAKE ANONYMITY FOR INEFFECTIVENESS



# PROFILING IS NOT ENOUGH

UNMASKING THESE FIVE "GHOSTS" IS A CRITICAL BREAKTHROUGH.  
BUT IT IS ONLY THE FIRST STEP.

RAW DATA



PREDICTIVE  
INTELLIGENCE



KNOWING WHAT THEY LOOK LIKE IS *REACTIVE*.  
PREDICTING WHAT THEY WILL DO NEXT IS *PROACTIVE*.



01

02

03

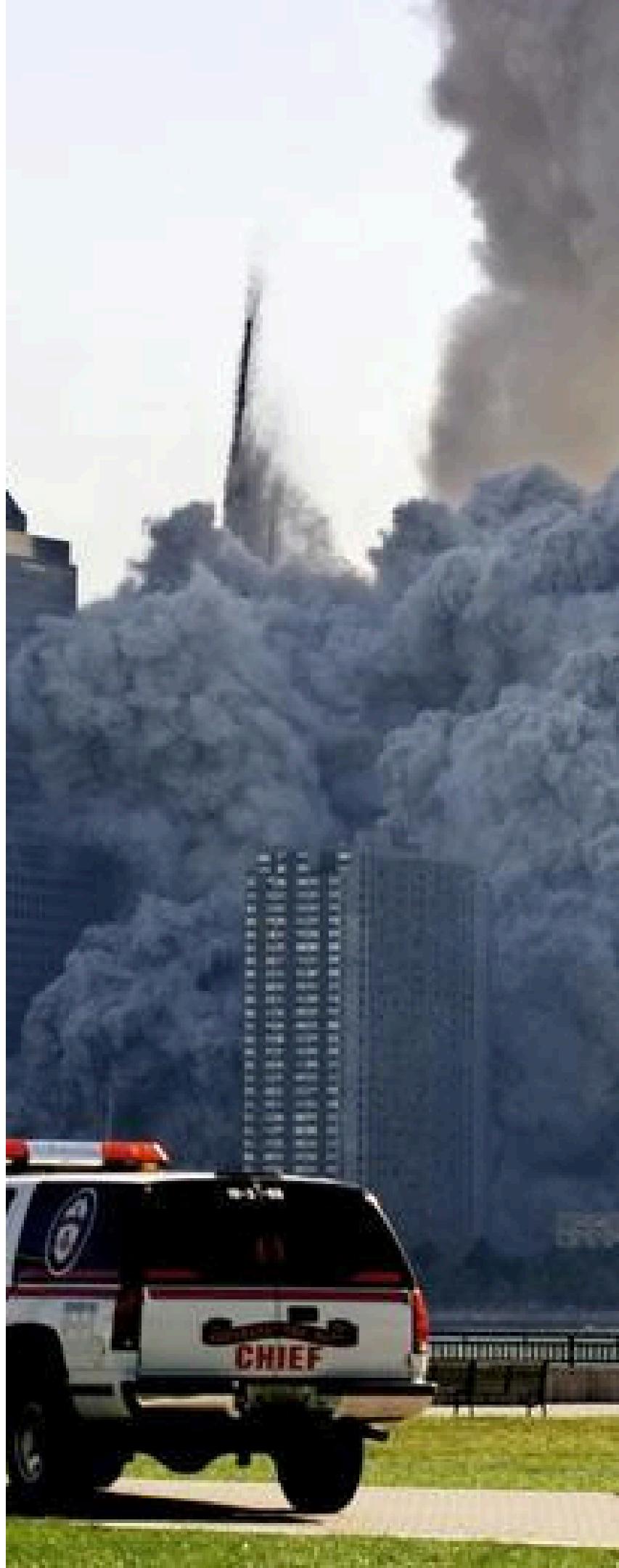
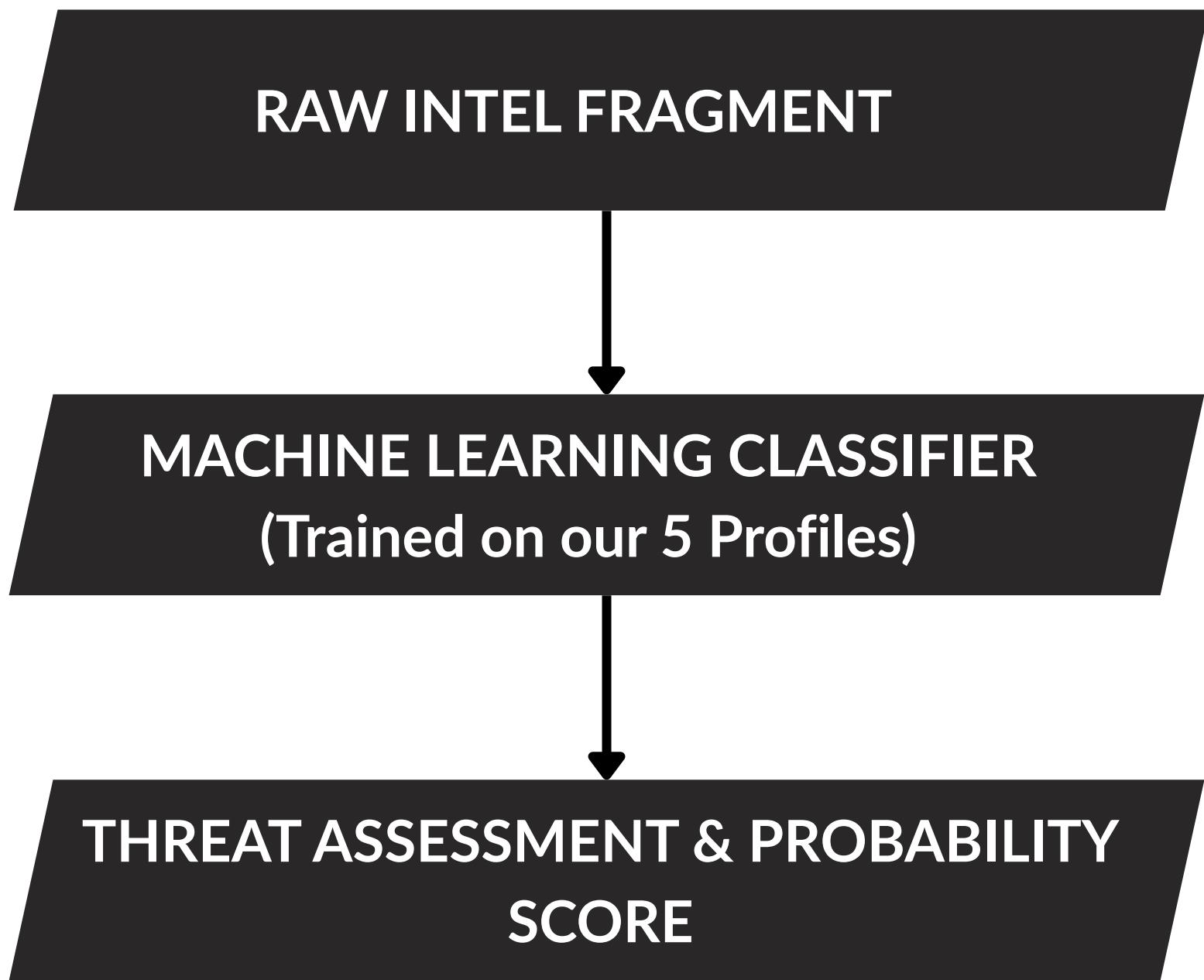
04

05

06

07

# THE BLUEPRINT: A PREDICTIVE EARLY WARNING SYSTEM



# THE FUTURE OF SECURITY IS NOT FOUGHT WITH *BULLETS*. IT IS WON WITH *DATA*.

SHIFT THE PARADIGM: FROM REACTIVE TO PREDICTIVE

BUILD THE SYSTEM: AI-POWERED EARLY WARNING SYSTEM

PRIORITIZE DATA PREPARATION: ACCURACY OF ANY PREDICTIVE  
SYSTEM IS ENTIRELY DEPENDENT ON THE DATA QUALITY