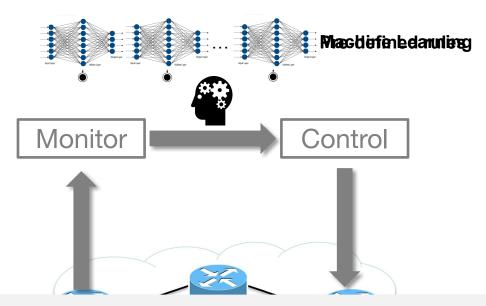
Making the Self-Driving "Net" Work

Developing Production-Ready ML Models for Self-Driving Networks

Arpit Gupta

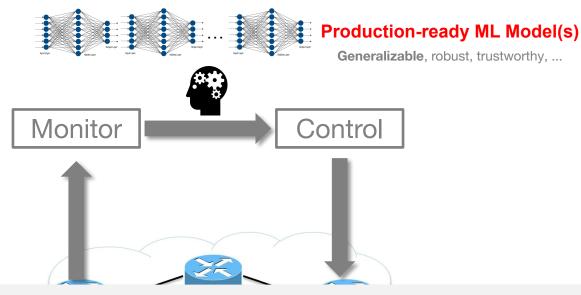
Associate Professor, UC Santa Barbara Faculty Scientist, Berkeley Lab (ESnet)

Self Driving Networks (AIOps)



Enable secure & performant connectivity with minimal human interventions

Key Requirement



Generalizable models perform as expected in target settings

AI/ML for Network Security: The Emperor has no Clothes

https://trusteeml.github.io/

Model Generalizability

Arthur S. Jacobs UFRGS, Brazil asjacobs@inf.ufrgs.br

Ronaldo A. Ferreira UFMS, Brazil raf@facom.ufms.br Roman Beltiukov UCSB, USA rbeltiukov@ucsb.edu

Arpit Gupta UCSB, USA arpitgupta@ucsb.edu Walter Willinger NIKSUN Inc., USA wwillinger@niksun.com

Lisandro Z. Granville UFRGS, Brazil granville@inf.ufrgs.br

Lack of generalizability attributable to at least 3 underspecification issues

Shortcut Learning

Model takes shortcuts to classify data (cheating)

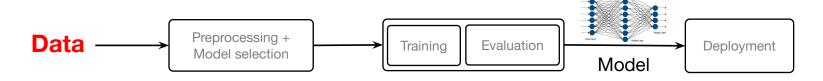
OOD Issues

Model does not generalize to out-of-distribution samples (rote learning)

ML models in networking are vulnerable to these underspecification issues!

Fundamental Roadblocks

Places the responsibility on users to find the **right data**, i.e., data that enables developing **generalizable** ML model

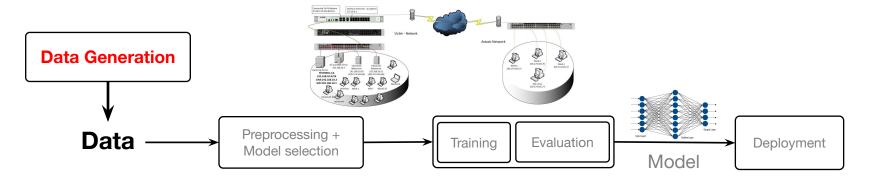


Standard ML Pipeline

For any given learning problem and target environment what is the "right" data

Fundamental Roadblocks

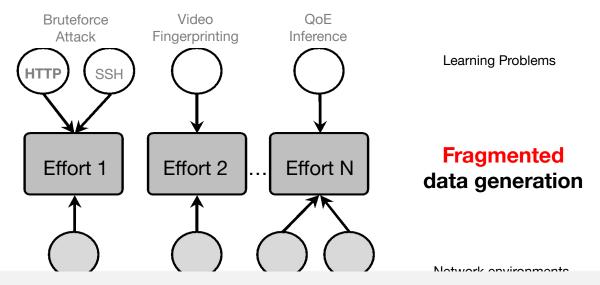
Networking problems necessitate **endogenous** data collection



Standard ML Pipeline

For any given learning problem and target environment how to generate the "right" data?

Data Generation Challenges

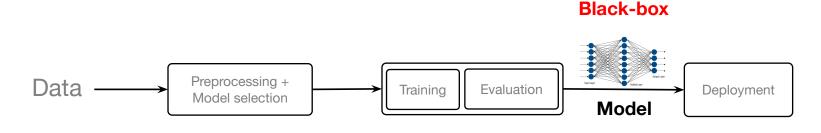


Overspecialized efforts with limited extensibility; public datasets dictate choice of problems & models

AWS Netrics UCSB

Fundamental Roadblocks

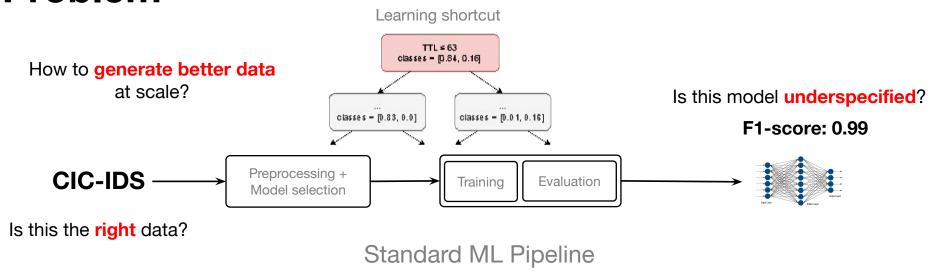
Outputs the most performant model, with little to no insights into model's decision-making (black-box)



Standard ML Pipeline

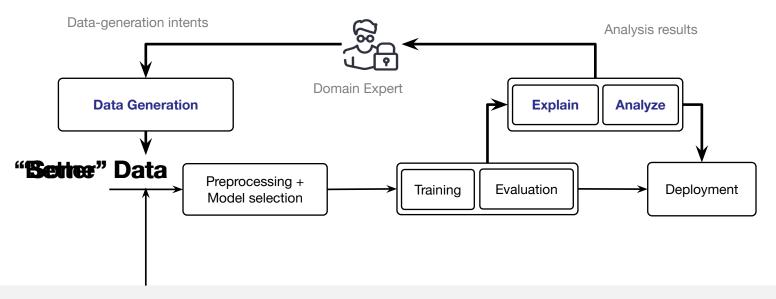
How to assess if resulting models are generalizable, i.e., not underspecificied?

Example: HTTP Brute Force Attack Detection Problem



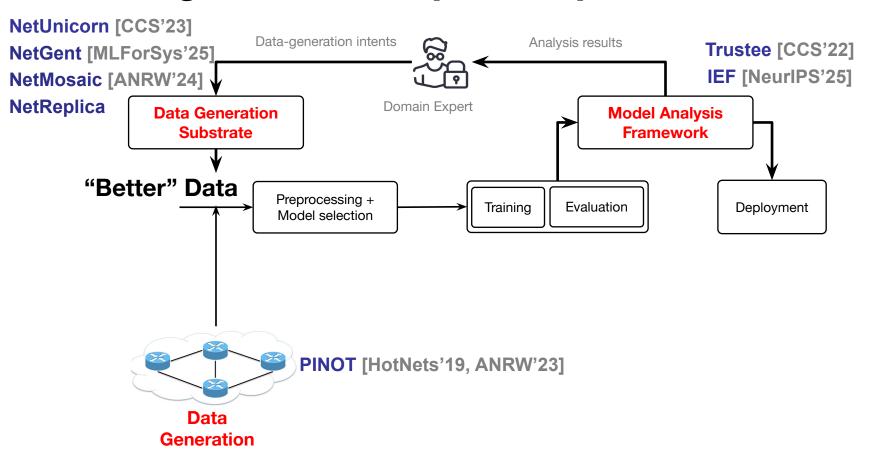
Answering these questions critical for developing generalizable ML artifacts for networking

Our Approach: Closed-Loop ML Pipeline



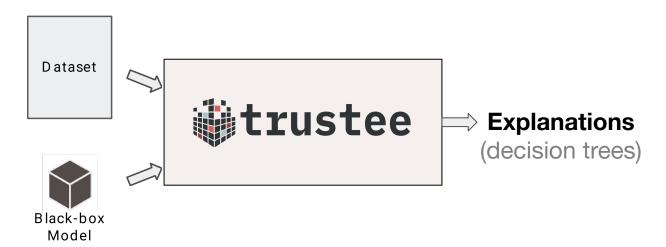
For any given learning problem and environment, iteratively fix underspecification issues to generate "better" data

Realizing Closed-Loop ML Pipeline



Trustee

A post-hoc global model explainability framework



Model-agnostic, high-fidelity, low-complexity, and stable decision trees to explain model's decision-making

Generates a Trust Report

								lassificat:	ion Trust	Report							
!	Summary																
i	Blackbox				i	Whitebox				Top-k Whitebox							
ı	Moi	del:	RandomFor	restClassif	ier	i	Explanation method: Trustee				Explanation method: Trustee						
	Datase ¹	t size:	9	947072			Mo	del:	Decision	TreeClassi	fier		Model: DecisionTreeClassifier			fier	
	Train/Te	st Split:	70.009	s / 30.00%			Iterations: 1			Iterations: 1							
							Sampl	e size:		50.00%			Samp1	le size:		50.00%	
							Decision	Tree Info					Decision	Tree Info			
							Si	ze:		2437				ize:			
					Depth: 31			Depth: 4									
					Leaves: 1219 # Input features: 18 (29.51%)			Leaves: 5									
	# Input features: 61 # Output classes: 5			Top-k: 1 # Input features: -													
							# Output classes: 5 (100.00%)										
											# Output	classes:		100.00%)			
i +	+ Performance			i +		E-	idelity		<u>:</u>	+			idelity				
¦			TOTHIGHEE		-				ļ								
		precision	recall	f1-score	support			precision	recall	f1-score	support			precision	recall	f1-score	support
ii	0	1.000	0.912	0.954	24408	i i	0	1.000	1.000	1.000	22254		0	0.000	0.000	0.000	22254
Π		0.752	0.910	0.824	1872	Ti		1.000	1.000	1.000	2265			0.000	0.000	0.000	2265
11		0.929	0.827	0.875	10994	11		0.969	0.965	0.967	9781			0.000	0.000	0.000	9781
Π		0.997	0.929	0.962	65188	i i		0.998	0.998	0.998	60768			0.544	0.957	0.694	60768
11		0.958	0.997	0.978	181660	i i		0.998	0.998	0.998	189054			0.875	0.821	0.847	189054
	accuracy			0.967	284122		accuracy			0.997	284122		accuracy			0.751	284122
1.1	macro avo	0.927	0.915	0.918	284122 I	1.1	macro avo	0.993	0.992	0.993	284122 I		macro avo	0.284	0.356	0.308	284122 I

Lowers the threshold of detecting underspecification issues

Trustee in Action

Problem	Dataset(s)	Model(s)	Issues
Detect VPN traffic	Public VPN dataset [20]	1-D CNN [61]	Shortcut learning
Detect Heartbleed traffic	CIC-IDS-2017 [54]	RF Classifier [54]	Out-of-distribution samples
Detect Malicious traffic (IDS)	CIC-IDS-2017 [54], Campus dataset	nPrintML [32]	Spurious correlations
Anomaly Detection	Mirai dataset [44]	Kitsune [44]	Out-of-distribution samples
OS Fingerprinting	CIC-IDS-2017 [54]	nPrintML [32]	Potential out-of-distribution samples
IoT Device Fingerprinting	UNSW-IoT [56]	Iisy [63]	Likely shortcut learning
Adaptive Bit-rate	HSDPA Norway [49]	Pensieve [42]	Potential out-of-distribution samples

Demonstrated prevalence of underspecification issues in existing ML artifacts for network security

I lodovoposification

Trustee in Action



IETF/IRTF Applied Networking Research Prize, 2023



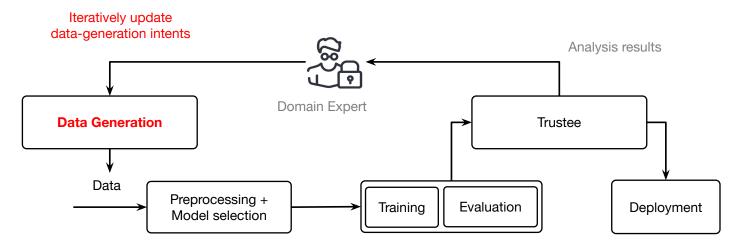
ACM CCS Best Paper Honorable Mention, 2022

Public Datasets

Problem	Dataset(s)	Model(s)	
Detect VPN traffic	Public VPN dataset [20]	1-D CNN [61]	Shortcut learning
Detect Heartbleed traffic	CIC-IDS-2017 [54]	RF Classifier [54]	Out-of-distribution samples
Detect Malicious traffic (IDS)	CIC-IDS-2017 [54], Campus dataset	nPrintML [32]	Spurious correlations
Anomaly Detection	Mirai dataset [44]	Kitsune [44]	Out-of-distribution samples
OS Fingerprinting	CIC-IDS-2017 [54]	nPrintML [32]	Potential out-of-distribution samples
IoT Device Fingerprinting	UNSW-IoT [56]	Iisy [63]	Likely shortcut learning
Adaptive Bit-rate	HSDPA Norway [49]	Pensieve [42]	Potential out-of-distribution samples

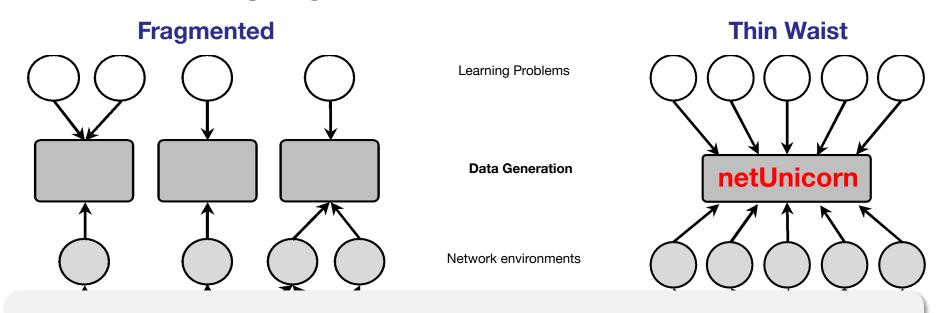
Attributed most underspecification issues to publicly available datasets

How to Fix the Underspecification (Data) Issues?



Use insights from model analysis to iteratively generate "better" data

How to Simplify Data Generation?



How to realize data-generation thin waist?

Disaggregation

Core Principle

A true thin waist emerges only when we decouple the layers that have historically been entangled

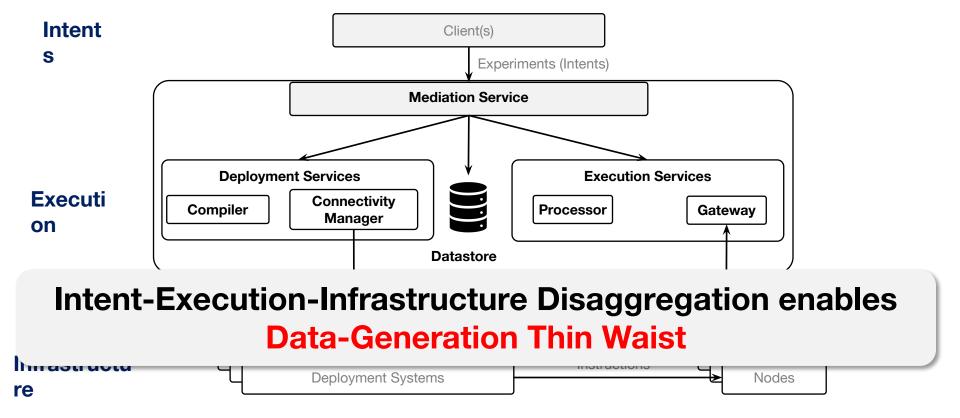
Disaggregate Intents from Execution

- Intent = what the experimenter wants: Execution = how it runs
- Decoupling makes intents portable, reusable, and stable across environments

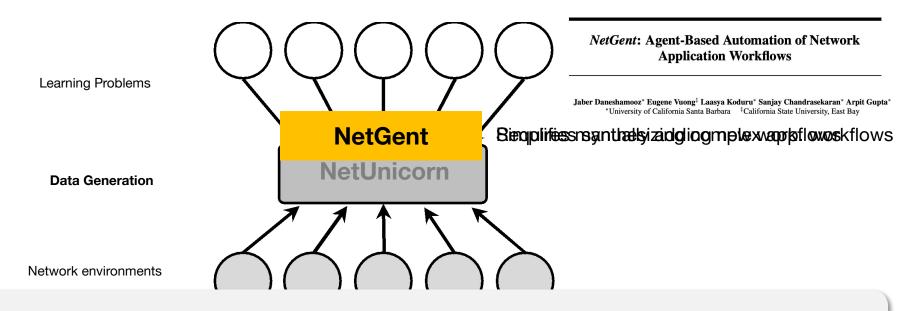
Disaggregate Execution from Infrastructure

Separation enables automatic mapping, capability-aware planning, and scaling across heterogeneity

Disaggregation In Action



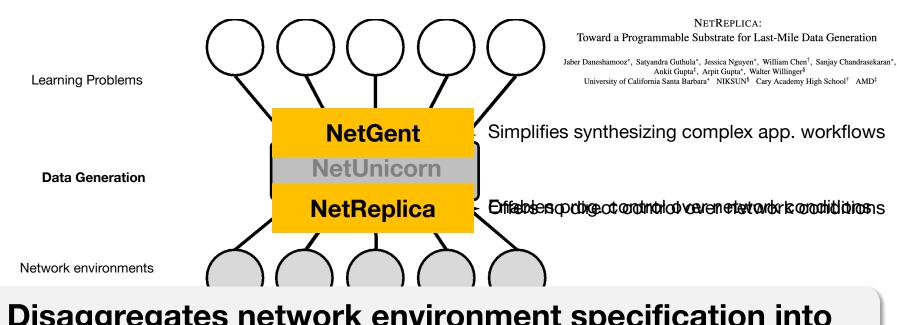
Augmenting Data Generation "Thin Waist"



Disaggregates application workflows into reusable tasks

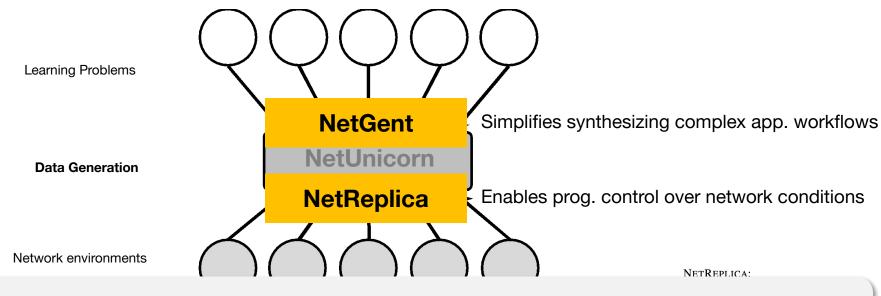


Augmenting Data Generation "Thin Waist"



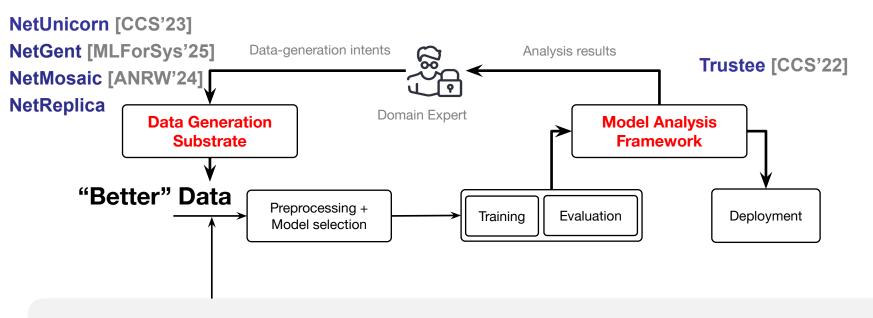
Disaggregates network environment specification into static and dynamic attributes

Augmenting Data Generation "Thin Waist"



Programmatically generate data for hundreds of applications across millions of different network conditions

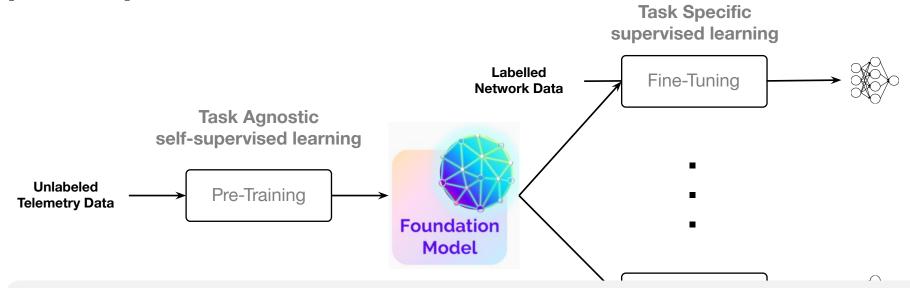
Closed-Loop ML Pipeline



Only fixes generalizability for one model at a time

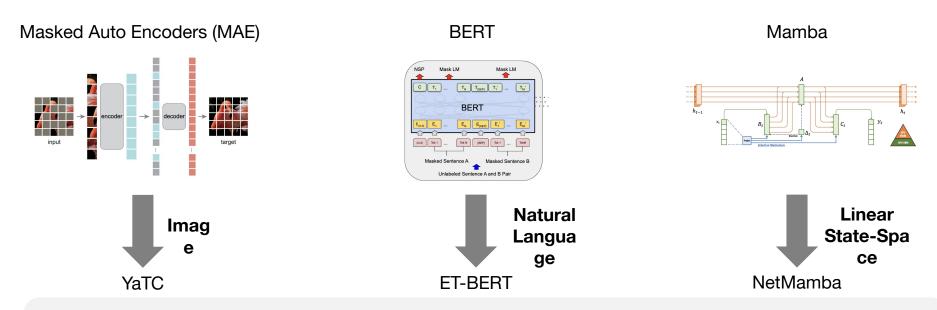
Data
Generation

New Frontier: Network Foundation Models (NFMs)



Critical to reason about generalizability of pre-trained Network Foundation Models

The Current Landscape for NFMs



Force mapping of network data to images, natural language, or linear state-space models





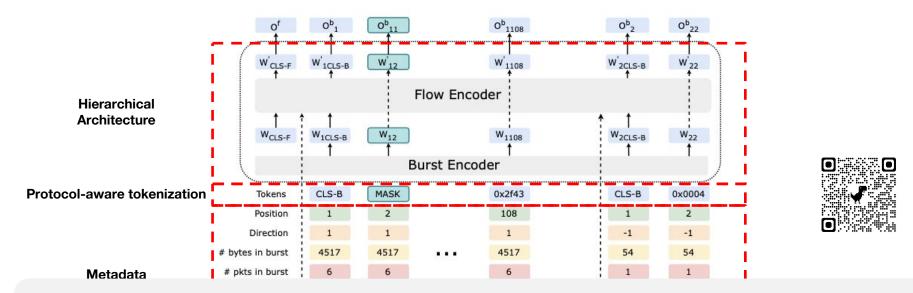


What's Unique about Network Data (Packet Traces)?

- Embody different protocols (tokenization?)
 - Packet content is dictated by disparate protocols and standards
- Entail variable length sequences (token selection?)
 - Packet sequences at any granularity is heavy tailed (multi-fractal behavior)
- Encode multi-modal information (token embedding?)
 - Packet sequences carry critical spatial, temporal, & contextual information
- Intrinsically hierarchical (modelling?)
 - Packets in different spatial/temporal groups have disparate semantic meanings

Necessitates a "Domain-specific" Approach

netFound: A Domain-Specific Network Foundation Model



How do we know if netFound is better, i.e., not only performant but also generalizable?

Standard Approach



Network

			N	Models	Foundatio	n Models		
Task	Туре	Dataset	Curtains (%)	nPrintML (%)	ET-BERT (%)	YaTC (%)	netFound (our) (%)	
1	Traffic Classification	Campus dataset	54.53 ± 0.97	87.22 ± 0.12	72.26 ± 0.38	76.54 ± 0.23	96.08 ± 0.04	
1	Tranic Classification	Campus dataset	p < 0.001	p < 0.001	p < 0.001	p < 0.001	_	
2		Crossmarkets [63] (Acc@10)	20.64 ± 0.13	64.83 ± 0.28	35.62 ± 0.39	58.13 ± 0.89	66.35 ± 0.99	
4	Application Fingerprinting	Crossmarkets [03] (Acc@10)	p < 0.001	p = 0.098	p < 0.001	p = 0.010	(- 3	
3	Application imgerprinting	ISCXVPN-2016 [18]	66.85 ± 2.21	84.10 ± 0.41	77.57 ± 1.20	83.84 ± 0.24	91.02 ± 0.10	
3			p = 0.003	p < 0.001	p < 0.001	p < 0.001	1.—1	
4	Intrusion Detection	CICIDS2017 [55]	99.75 ± 0.16	99.93 ± 0.01	99.94 ± 0.01	99.92 ± 0.01	99.99 ± 0.01	
4	intrusion Detection	CICID32017 [55]	p = 0.082	p = 0.012	p = 0.018	p = 0.005	-	
5	HTTP Bruteforce Detection	netUnicorn [11]	96.82 ± 0.22	98.51 ± 0.02	98.63 ± 0.02	98.73 ± 0.10	99.01 ± 0.01	
3	III IF Bluteloice Detection	netoincom [11]	p = 0.006	p < 0.001	p < 0.001	p = 0.030	=	
		2	•			~		

Deep Learning

Exhibits better performance for diverse "well-explored" learning tasks and datasets

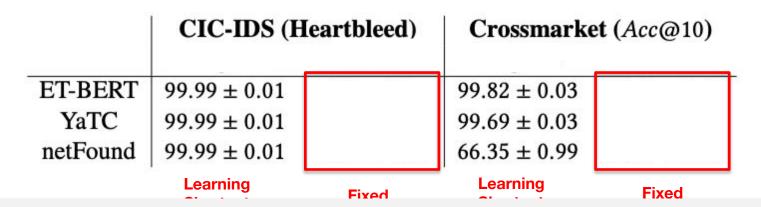
Wait! Which Task? What Dat

The Sweet Danger of Sugar: Debunking Representation Learning for Encrypted Traffic Classification

Yuqi Zhao Politecnico di Torino Torino, Italy yuqi.zhao@polito.it Giovanni Dettori Politecnico di Torino Torino, Italy giovanni.dettori@polito.it Matteo Boffa Politecnico di Torino Torino, Italy matteo.boffa@polito.it

Luca Vassio
Politecnico di Torino
Torino, Italy
luca.vassio@polito.it

Marco Mellia Politecnico di Torino Torino, Italy marco.mellia@polito.it



Equating NFMs' success with performance on a limited set of downstream tasks and datasets is misleading

Our Approach: Intrinsic Evaluation Framework (IEF)

Assess embedding quality decoupled from downstream tasks/datasets

Embedding Geometry Analysis

How efficiently the model utilizes representation space?

Metric Alignment Assessment

How well it captures well-known features, developed by experts over time?

Causal Sensitivity Testing

How sensitive is the model to network context perturbations?

Demystifying Network Foundation Models

Experimental Setup

Sylee (Roman) Beltiukov* UC Santa Barbara Satyandra Guthula UC Santa Barbara

Wenbo Guo UC Santa Barbara

Walter Willinger NIKSUN, Inc Arpit Gupta UC Santa Barbara

Network Foundation Models

YaTC, ET-BERT, NetFound, and NetMamba



Datasets

Endogenous (Actively Generated)

Android Crossmarket, CIC-IDS, APT-IIoT

Exogenous (Passive Traces)

MAWI, CAIDA

Embedding Geometry Analysis

Claim

Effective NFMs' embeddings should utilize the full representation space

Rationale

Well-distributed embedding geometry indicates that model distinguishes between flow along meaningful dimensions

Embedding and Cosine Similarity Embedding Vectors Cat Dog Cosine Similarity

Methodology

Measure anisotropy (cosine similarity) between embeddings pairs in the dataset

Example of embedding distribution in the space: *less related embeddings* are further from each other

Embedding Geometry Analysis Results

	Yal	ГС	ET-BERT	netFound	NetMamba
Dataset	cos	•	cos	cos	cos
Crossmarket	0.85		0.88	0.69	0.93
CIC-APT-IIoT24	0.87		0.88	0.82	0.98
CIC-IDS2017	0.85		0.74	0.69	0.92
CAIDA	0.87		0.71	0.86	0.99
MAWI	0.88		0.78	0.94	0.99

We observe total representation collapse (NetMamba), dimensional dominance (YaTC), and limited generalizability (netFound)

Metric Alignment Assessment

Claim

Effective NFMs' should encode critical network statistics such as flow duration, packet size distributions, TCP dynamics, etc., without explicit supervision.

Rationale

These features, developed over decades, encode domain-knowledge, which worked well for different learning tasks

PCAP CICFlowme Foundation Model **Embeddings Features** Correlation (CKA)

Methodology

Compute **CKA** similarity between each established metric and embeddings across flows

Metric Alignment Assessment

	Crossmarket	CIC-APT-IIoT24	CIC-IDS2017
YaTC	0.098	0.148	0.092
ET-BERT	0.012	0.014	0.064
netFound	0.156	0.219	0.167
NetMamba	0.047	0.141	0.042

Except netFound, most NFMs miss established metrics; yet netFound falters on production traces.

Result Summary: Intrinsic Evaluation Framework

Embedding Geometry Analysis

Total collapse (NetMamba), dimensional dominance (YaTC), limited generalizability (netFound)

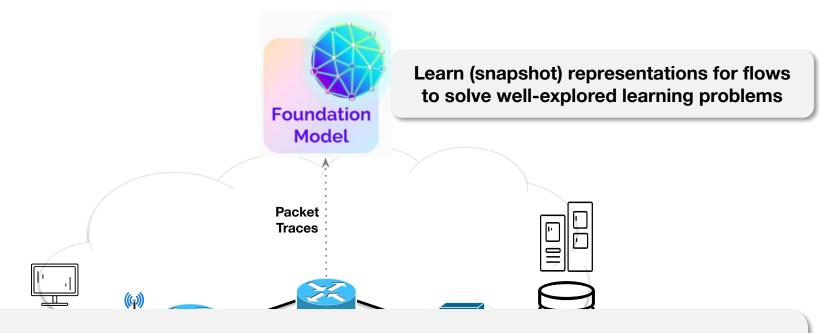
Metric Alignment Assessment

Poor alignment for most except netFound, which struggles to generalize

Causal Sansitivity Testing

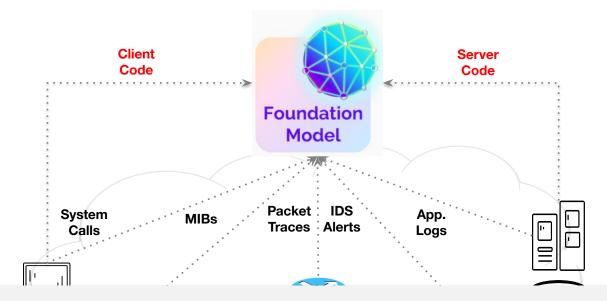
While no clear winner emerges, this framework provides an insightful benchmark

Current State



Can we learn more expressive and generalizable representations with least complexity and minimal data?

Where can we go from here?



Learn more expressive representations to facilitate solving "unexplored" learning problems in networking



"Unexplored" Learning Problems in Networking

Less Ambitious

Facilitate interactions with network data in natural language---rethink network telemetry systems

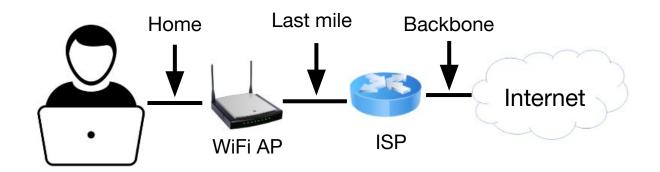
Moderately Ambitious

Identify and diagnose subtle and complex network events (anomalies), predict future events, and search/rank similar past events

Caution: Please remember the "garbage-in-garbage-out" adage and don't blindly throw LLMs for these problems

Code-2-pkt & pkt-2-code transformations for verification

Exemplary Use Case

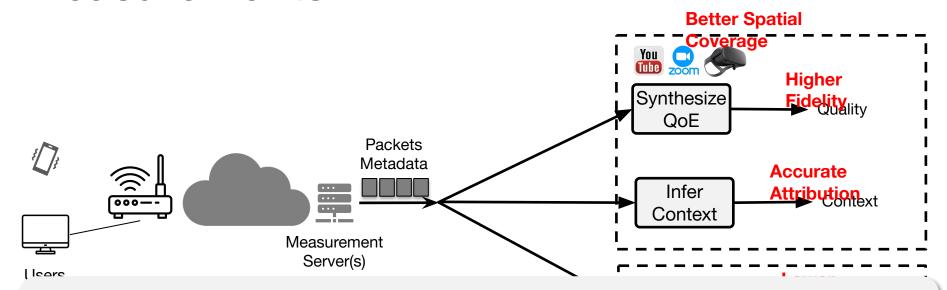


For every residence (i.e., **dense spatial coverage**), periodically (i.e., **dense temporal coverage**)

Report metric(s) that captures users' quality of experience* (high-fidelity assessment)

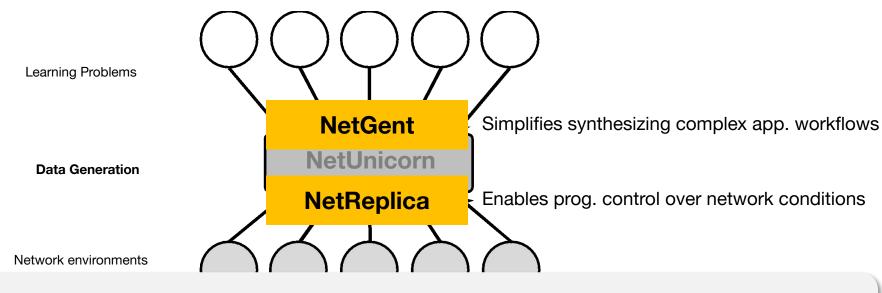
Existing tools fail to bridge gap between measured and experienced broadband quality

Rethinking Broadband Quality Measurements



Learn data representations that enable closed-loop measurements, context inference, and counterfactual analysis

Programmable Data Generation Substrate



Provides a unique tool to pursue this "grand challenge" problem in network measurement



Summary



NeurIPS' Demystifying Network Foundation Models

Sylee (Roman) Beltiukov* UC Santa Barbara Satyandra Guthula UC Santa Barbara

Wenbo Guo UC Santa Barbara Walter Willinger NIKSUN, Inc Arpit Gupta UC Santa Barbara

- Production-ready ML models critical for self-driving networks
- Closed-loop ML pipeline, composed of model-analysis (Trustee) and data-generation (NetUnicorn, NetReplica, NetGent) substrates critical for production-ready ML
- Network foundation model(s) are critical for self-driving networks
 - Intrinsic evaluation framework helps answer what a Network foundation model is (and is not) learning
 - Programmable data-generation substrate provide a unique opportunity to take on new and revisit older problems with a fresher perspective