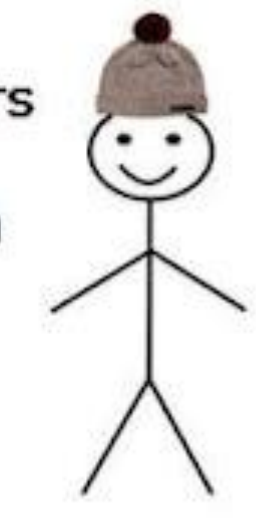


This is Reviewer 2.
 Reviewer 2 is an angry and bitter scholar exacting revenge on their peers through overly critical anonymous rejections of papers they secretly wish they would have written.
 Reviewer 2 does not like puppies.
 Don't be like Reviewer 2.



When Adversarial Perturbations meet Concept Drift: An Exploratory Analysis on ML-NIDS

Giovanni Apruzzese*, Aurore Fass§, Fabio Pierazzi||

*University of Liechtenstein, §CISPA Helmholtz Center of Information Security, ||King's College London

giovanni.apruzzese@uni.li*, fass@cispa.de§, fabio.pierazzi@kcl.ac.uk||

Shoutout to our Reviewer 2, who wrote: "From my point of view, if all ML-NIDS papers were as well documented as this one, research in ML-NIDS would gain a great deal!"

Review #108
 Overall merit
 5. Strong accept

Our Institutions



Abstract

We scrutinize the effects of "blind" adversarial perturbations against machine learning (ML)-based network intrusion detection systems (NIDS) affected by concept drift. There may be cases in which a real attacker – unable to access and hence unaware that the ML-NIDS is weakened by concept drift – attempts to evade the ML-NIDS with data perturbations. It is currently unknown if the cumulative effect of such adversarial perturbations and concept drift leads to a greater or lower impact on ML-NIDS. In this "open problem" paper, we seek to investigate this unusual, but realistic, setting—we are not interested in perfect knowledge attackers.

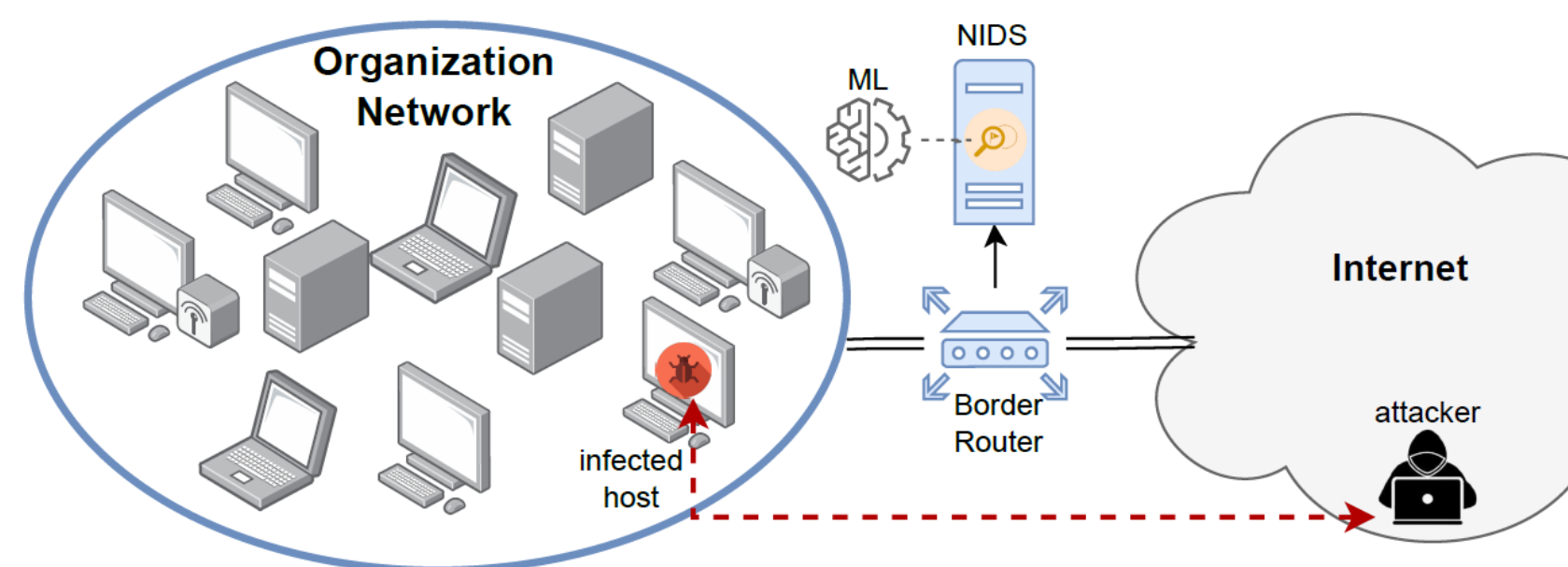
We begin by retrieving a publicly available dataset of documented network traces captured in a real, large (>300 hosts) organization. Overall, these traces include several years of raw traffic packets—both benign and malicious. Then, we adversarially manipulate malicious packets with problem-space perturbations, representing a physically realizable attack. Finally, we carry out the first exploratory analysis focused on comparing the effects of our "adversarial examples" with their respective unperturbed malicious variants in concept-drift scenarios. Through two case studies (a "short-term" one of 8 days; and a "long-term" one of 4 years) encompassing 48 detector variants, we find that, although our perturbations induce a lower detection rate in concept-drift scenarios, some perturbations yield adverse effects for the attacker in intriguing use cases. Overall our study shows that the topics covered are still an open problem which require a re-assessment from future research.

What do we do (and why do we do it)?

Contributions

- We shed light on a problem that has never been investigated before in the ML-NIDS context: the combination of realistic (blind, realizable) adversarial perturbations with concept drift. We:
- pinpoint an open-source (and documented) dataset that can be used for concept-drift assessments in ML-NIDS (and which has been overlooked by most research);
- craft problem-space adversarial perturbations by manipulating raw network traffic simulating a simple and feasible attack;
- investigate the extent to which ML-NIDS are statistically significantly affected by realistic adversarial perturbations in concept drift contexts

Threat Model



The attacker is outside the organization's network, cannot access the NIDS (e.g., to query it), and only knows that an ML model analyzes network traffic. The attacker knows that the ML-NIDS is trained on datapoints "similar" to those used in their malicious activities; hence, the attacker is aware that there is a risk of being detected if nothing is done. However, the attacker does not know if the ML-NIDS is affected (or not) by concept drift (and neither do the defenders).

Motivation

A real-attacker's dilemma. It is unknown whether it is beneficial to introduce blind adversarial perturbations in an attempt to bypass an ML-NIDS affected by concept drift. The attacker must make a choice—potentially a detrimental one!

Challenge: Finding the right Data

How to carry out meaningful concept-drift assessments?

Concept drift pertains to supervised ML, and defines the temporal "degradation" of an ML-model's performance due to naturally occurring phenomena. Hence, an assessment requires:

- Data that captures a "relatively long" timeframe
- Data that reflects "naturally occurring" phenomena
- Data that is "unambiguously" labeled.

Unfortunately, meeting such a threefold requirement is tough in the ML-NIDS context.

Data "pitfalls": what can existing datasets for ML-NIDS allow?

We surveyed the landscape of existing publicly available datasets for ML-NIDS, and we found that most do not allow a "meaningful" concept drift assessment. Some examples:

- NSL-KDD is flawed ☹
- CICIDS17 is captured over a timespan of 5 days (!), and the data was created via simulations
- UNSW-NB15 is collected over 15 hours (!)
- Kitsune also has barely 1 day of data.
- UGR16 is captured over 100 days...but the labeling is not consistent and is not provided with PCAP intriguingly, most prior research on concept drift was evaluated on the abovementioned datasets!

The MCFP data: a practical solution (for our community)

We found that the best publicly-available dataset that meets all our requirements is the "Malware Capture Facility Project" (<https://www.stratosphereips.org/datasets-malware>), which is an extension of the popular CTU13 dataset.

Takeaway. The MCFP is a suitable solution for assessments of ML-NIDS under concept drift. It is large, entails long timespans, ground truth is provided, and includes various types of benign (from hundreds of hosts) and malicious (entailing recent attacks) traffic in PCAP format – hence useful for "problem-space" attacks.

Perhaps surprisingly, we found only one work ([49]) on "concept drift" that considers MCFP (but it is unclear how the temporal aspect was taken into account).

A "snippet" of the MCFP dataset (which we will use for our analysis)

Trace (Link)	Date	PCAP Size	Flows	Nature
1	17 Dec 2013	400M	10K	Background
2	17 Dec 2013	800M	10K	Active
3	24 Mar 2015	1M	1K	Background
4	13 Sep 2016	63M	40	Background
5	13 Sep 2016	63M	40	Active
6	13 Sep 2016	0.5M	40	Background
7	18 Apr 2017	200M	10K	Background
8	19 Apr 2017	200M	7.8K	Active
9	25 Apr 2017	400M	7.8K	Background
10	26 Apr 2017	110M	10K	Active
11	30 Apr 2017	270M	21K	Background
12	30 Apr 2017	424M	39K	Active
13	1 May 2017	173M	12K	Background
14	1 May 2017	435M	35K	Active
15	1 May 2017	800M	55K	Background
16	1 May 2017	725M	60K	Active
17	2 May 2017	300M	15K	Background
18	2 May 2017	116M	105K	Active
19	2 May 2017	820M	82K	Background
20	3 Jul 2017	62M	24	Active
21	23 Jul 2017	400M	62K	Background
22	5 Sept 2017	16M	1K	Active
23	7 May 2018	800M	43K	Background

Table 6: Malicious PCAP traces from MCFP used in our assessment (we will perturb these). Note: the most recent traces in MCFP are collected in 2021. We provide an explanation of the method used to choose these traces in our supplementary document [3].

m	Trace (Link)	Date	PCAP Size	Flows	Total		
1	14 May 2017	500K	5K	37M	23K		
2	14 May 2017	11M	15K	336M	30K		
3	15 May 2017	3.6M	171	772M	226K		
4	15 May 2017	11M	32K	16 Aug 2017	153M	11K	
5	24 May 2017	444M	9K	16 Aug 2017	146M	10K	
6	11 Jul 2017	1.6M	14K	1	24 Jun 2017	52M	24K
7	11 Jul 2017	7.6M	17K	2	3 Aug 2017	6.4M	2K
8	11 Jul 2017	7.3M	13K	2	3 Aug 2017	252M	69K
9	11 Jul 2017	4.8M	93K	1	29 Jun 2017	83M	40K
10	11 Jul 2017	3.1M	35	2	30 Mar 2017	90M	41K
11	11 Jul 2017	6.3M	4K	3	30 Mar 2017	90M	41K
12	11 Jul 2017	6.4M	17K	13	24 Jun 2017	74M	33K
13	11 Jul 2017	1.6M	17K	4	12 Apr 2017	288M	160K
14	12 Jul 2017	6.1M	3.6K	6	12 Apr 2017	115M	53K
15	13 Jul 2017	6.2M	210	6	17 Apr 2017	142M	103K
16	13 Jul 2017	6.8M	11K	7	17 Apr 2017	142M	103K
17	13 Jul 2017	6.7M	10K	8	8 May 2017	214M	127K
18	13 Feb 2017	79M	102	9	15 May 2017	204M	79K
19	11 Aug 2017	57M	25K	10	7 Jun 2017	21M	124K
20	11 Aug 2017	31M	51K	11	15 Jun 2017	225M	141K
21	18 Apr 2017	66M	30K	12	24 Jun 2017	77M	31K
22	18 Apr 2017	47M	35K	13	24 Jun 2017	74M	33K
23	15 May 2017	7.4M	43K	14	24 Jun 2017	78M	31K
24	15 May 2017	33M	48K	15	24 Jun 2017	44M	27K
25	16 May 2017	52M	63K	16	30 Jun 2018	33M	19K
26	24 Jun 2017	16M	11K	17	30 Jun 2018	212M	62K
27	29 Jun 2018	310M	73K	18	2 Feb 2018	197M	59K
28	30 Jun 2018	190M	37K	19	27 Mar 2018	410M	122K
29	30 Jun 2018	223M	52K	20	30 Jul 2021	220M	41

m	Trace (Link)	Date	PCAP Size	Flows	Total
1	10 Aug 2011	6.1G	4M	35K	Background
2	10 Aug 2011	6.2G	2.5M	10K	Active
3	11 Aug 2011	6.2G	5M	11K	Background
4	12 Aug 2011	14.5G	5M	113K	Background
5	15 Aug 2011	5.4G	1.4M	30K	Active
6	15 Aug 2011	0.4G	150K	5K	Background
7	16 Aug 2011	4.4G	2.1M	28K	Active
8	16 Aug 2011	9.8G	2.6M	35K	Active
9	17 Aug 2011	9.8G	2.6M	35K	Active
10	18 Aug 2011	0.6G	113K	2.8K	Background

Realistic Perturbations and Case Studies

Crafting "blind" and "realizable" adversarial perturbations

We only manipulate malicious packets (i.e., those that can be detected by the ML-NIDS). For each PCAP trace we consider, we create four "adversarial" traces by proceeding as follows:

- Take all UDP packets and append a small padding of [1-100] (random) bytes
- Take all TCP packets with the PSH flag active, and apply the [1-100] bytes random padding
- Repeat the process again to mitigate bias due to randomness.

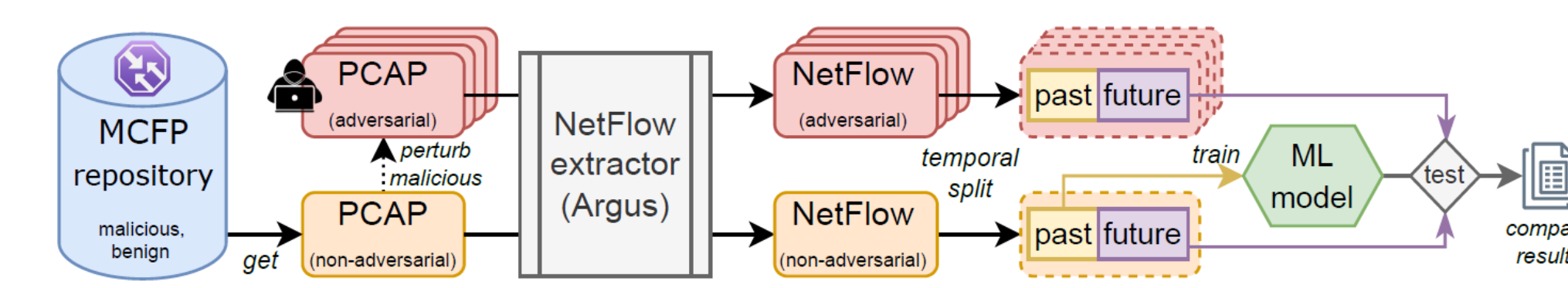
We also ensure each packet does not exceed the maximum packet length and recreate the checksum. We do these operations via *scapy* and we release our "packet manipulator" tool.

We consider ML-NIDS analysing NetFlows. Hence, the perturbations above will translate in the "feature space". We generate the NetFlows by using *Argus*, and we then label them by following the official documentation provided by MCFP's creators.

Moreover, in our "adversarial-evaluation" analyses, we will filter to only consider traffic generated from the attacker-controlled machine

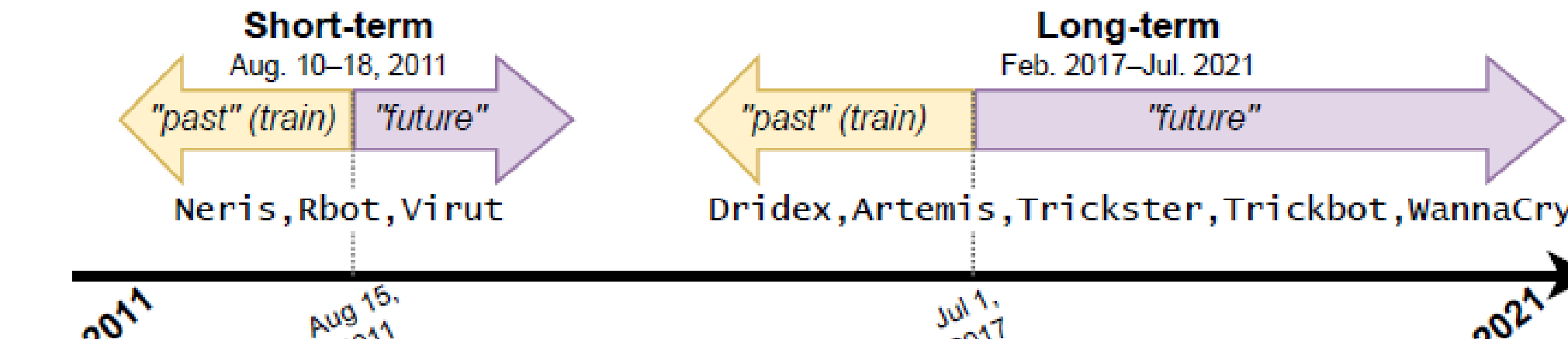
Workflow and ML-NIDS development

We do a "temporal split": the data before a certain date is considered "past data" and is used to develop the ML models used in the ML-NIDS. The data after a certain date is considered "future data" and is used to test our hypotheses during our analysis. Our models are binary classifiers (benign/malicious), and we consider 48 different ML-NIDS, varying the binary classifier (Random Forest of Histogram Gradient Boosting) and detector architecture (a single binary classifier, or an ensemble of classifiers—each devoted to a specific malware variant). We also consider an attack-agnostic defense [7] proven to work against adversarial ML attacks in the feature space.



Summary. We get benign and malicious PCAP from MCFP and manipulate only the malicious traces—yielding problem-space adversarial perturbations. Next, we take all PCAPs and extract (and label) the corresponding NetFlows. Finally, we train ML models (on both benign and malicious "past" data) and test their effectiveness on "future" data (benign, malicious, and adversarial) in a concept-drift setting. We repeat our experiments 50x to ensure statistically sound comparisons.

Overview of our analysis



Case Studies. We consider two case studies, capturing different time periods:

- Short-term Case Study (SCS), of one week (Aug. 10→18, 2011).
- The cutoff date (past/future data) is set to Aug. 15th 2011
- Long-term Case Study (LCS), of four years (Feb. 2017→Jul 2021).
- The cutoff date (past/future data) is set to Jul. 1st 2017

For SCS, we consider three malware types (Neris, Rbot, Virut). For LCS, we consider five (Trickster, Trickbot, Artemis, Wannacy, Dridex).

Preliminary Analysis: are our ML-NIDS good and is there Concept Drift?

Table 1: Pre-deployment results. We assess the performance of our ML-NIDS on the test set from "past" data. We report the average *tp* (on malicious samples) and *rr* (on benign samples). Cells in boldface are more relevant (they represent architectures denoting "generic" detectors). The defense is denoted with a \square .

		Full	Ens	Artemis	Dridex	Trickster	Trickbot	Wannacy
Benign	RF	0.999	0.999	0.999	1.000	1.000	0.999	0.999
Benign	HGB	0.999	0.999	0.999	0.999	0.999	0.999	0.999
Malicious	RF	0.994	0.992	0.993	0.998	0.964	0.999	0.999
Malicious	HGB	0.992	0.982	0.995	0.999	0.763	0.999	0.999
Malicious	CRF	0.999	0.999	0.999	0.999	0.999	0.999	0.999
Malicious	CHGB	0.999	0.999	0.999	0.999	0.999	0.999	0.999
Malicious	CRF	0.993	0.992	0.992	0.953	0.999	0.999	0.999
Malicious	CHGB	0.989	0.983	0.992	0.968	0.762	0.999	0.999

Table 2: Concept-drift results. We assess the performance of our ML-NIDS on the test set from "future" data. We report the *tp* (malicious) and *rr* (benign) averaged over 50 trials. Cells in red report cases in which the performance is statistically significantly ($p < 0.05$) worse than on the test set of "past" data (in Table 1). The defense is denoted with a \square .

		Full	Ens	Artemis	Dridex	Trickster	Trickbot	Wannacy
Benign	RF	0.989	0.993	0.993	1.000	1.000	0.986	0.999
Benign	HGB	0.990	0.982	0.989	0.999	0.999	0.981	0.993
Malicious	RF	0.675	0.567	0.701	0.828	0.693	0.827	0.888
Malicious	HGB	0.673	0.577	0.768	0.820	0.663	0.859	0.911
Malicious	CRF	0.990	0.996	0.996	0.999	0.999	0.955	0.965
Malicious	CHGB	0.983	0.991	0.995	0.998	0.995	0.956	0.962
Malicious	CRF	0.631	0.488	0.382	0.624	0.769	0.786	0.937
Malicious	CHGB	0.634	0.561	0.392	0.625	0.665	0.793	0.938

Answer: yes, our ML-NIDS are good; and yes, there is concept drift in both case studies!

Major Findings and Recommendations

Do adv. perturbations help against ML-NIDS affected by concept drift?

Table 3: Non-adversarial results. We measure the average *tp* on the "future" malicious NetFlows. We only consider UDP/TCP NetFlows starting from within the network. The defense is denoted with a \square . A \square denotes when the defense is statistically significantly ($p < 0.05$) better/worse than the "vanilla".

		Full	Ens	Artemis	Dridex	Trickster	Trickbot	Wannacy
UDP	RF	0.980	0.791	0.941	0.151	0.384	0.823	0.807
UDP	HGB	0.824	0.816	0.990	0.000	0.039	0.754	0.893
TCP	RF	0.929	0.715	0.366	0.436	0.773	0.919	0.988
TCP	HGB	0.943	0.863	0.858	0.625	0.861	0.983	0.983
UDP	CRF	0.902	0.394	0.011	0.117	0.000	0.696	0.981
UDP	CHGB	0.921	0.478	0.025	0.103	0.117	0.795	0.999
TCP	CRF	0.880	0.811	0.412	0.513	0.951	0.979	0.958
TCP	CHGB	0.867	0.810	0.374	0.507	0.858	0.982	0.976

Table 4: Adversarial results. We compute the average *tp* on "future" adversarial NetFlows. Cells in red denote cases in which the *tp* is statistically significantly worse than the baseline in Table 3. A \square denotes when the defense is statistically significantly ($p < 0.05$) better/worse than the "vanilla".

		Full	Ens	Artemis	Dridex	Trickster	Trickbot	Wannacy
UDP	RF	0.921	0.602	0.879	0.011	0.032	0.812	0.790
UDP	HGB	0.824	0.803	0.975	0.014	0.059	0.740	0.899
TCP	RF	0.902	0.717	0.381	0.292	0.769	0.929	0.934
TCP	HGB	0.934	0.863	0.472	0.470	0.822	0.967	0.980
UDP	CRF	0.873	0.344	0.022	0.050	0.000	0.698	0.981
UDP	CHGB	0.889	0.389	0.060	0.077	0.058	0.707	0.999
TCP	CRF	0.868	0.834	0.441	0.410	0.931	0.974	0.931