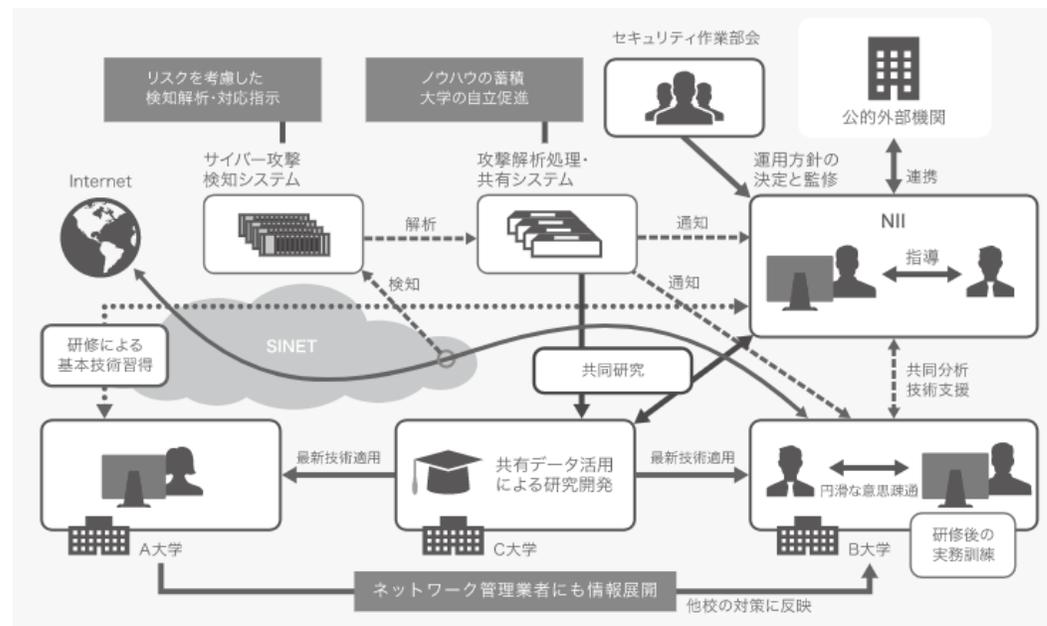


NII-SOCSベンチマークデータの 契約から研究発表まで

名古屋大学 情報基盤センター
情報基盤ネットワーク研究部門
基盤ネットワーク研究グループ
嶋田 創

NII-SOCSの設置目的に含まれる 「共有データ活用による研究開発」(1/3)

- 大学間が連携するための環境整備及び参加機関が学内のサイバーセキュリティ体制を確立するための支援[1]
 - 各大学担当者によるNII-SOCS上の自機関データの分析
 - NII-SOCSチームによる大学側担当者の教育(研修)や対応補助
 - **共有データ活用による研究開発**



[1]より抜粋した図

NII-SOCSの設置目的に含まれる 「共有データ活用による研究開発」(2/3)

- 設立当初の資料[1]でも、研究用データの公開はうたわれている

● 研究用データの公開

- 統計化・匿名化処理を施したベンチマークデータ
- バラマキ型の新種マルウェア情報の大学への提供
 - ・ 情報不足への対応も

[1]のスライド17より

● 警報情報とセッション情報

- 一般公開
 - ・ IPアドレス無しの単なる統計データ
- NDAに基づく研究機関向け公開
 - ・ IPアドレスはサニタイズ
 - ・ 観測時間を意図的に変動
 - ・ 警報の「通信の内容」部分は暗号化ままでハッシュ値を生成
 - ・ KyotoData2006+準拠

IPアドレスのサニタイズ
Saltを毎月変更
/24での連続性は保証

攻撃者自身による特定作業を防止

[1]のスライド19より

[1] 大学間連携に基づく情報セキュリティ体制の基盤構築 (H29 SINET・学術情報基盤サービス説明会)

https://www.nii.ac.jp/service/upload/7_setsumeikai2017_security_20171114.pdf

NII-SOCSの設置目的に含まれる 「共有データ活用による研究開発」(3/3)

- 研究用データは2021年より提供開始[1]
 - 匿名化処理をしたセッション情報+警報情報
 - 警報情報には、「x日後にパターンファイル更新して再度検知処理したらひっかかった」ものをゼロデイ攻撃としてカテゴリイズしたものもある
- 約款も準備されている

約款

検知情報に関する通知機能関連

- Web-API 機能利用約款  (令和2年11月30日制定)
- マルウェア警報情報ダウンロード機能利用約款  (令和2年11月30日制定)

研究用データ関連

- ベンチマークデータダウンロードサイト利用約款  (令和2年11月30日制定)
- マルウェア情報ダウンロードサイト利用約款  (令和2年11月30日制定)

[1] 国立情報学研究所、NII-SOCS の蓄積データを研究者向けに提供
<https://scan.netsecurity.ne.jp/article/2021/01/21/45077.html>

概要



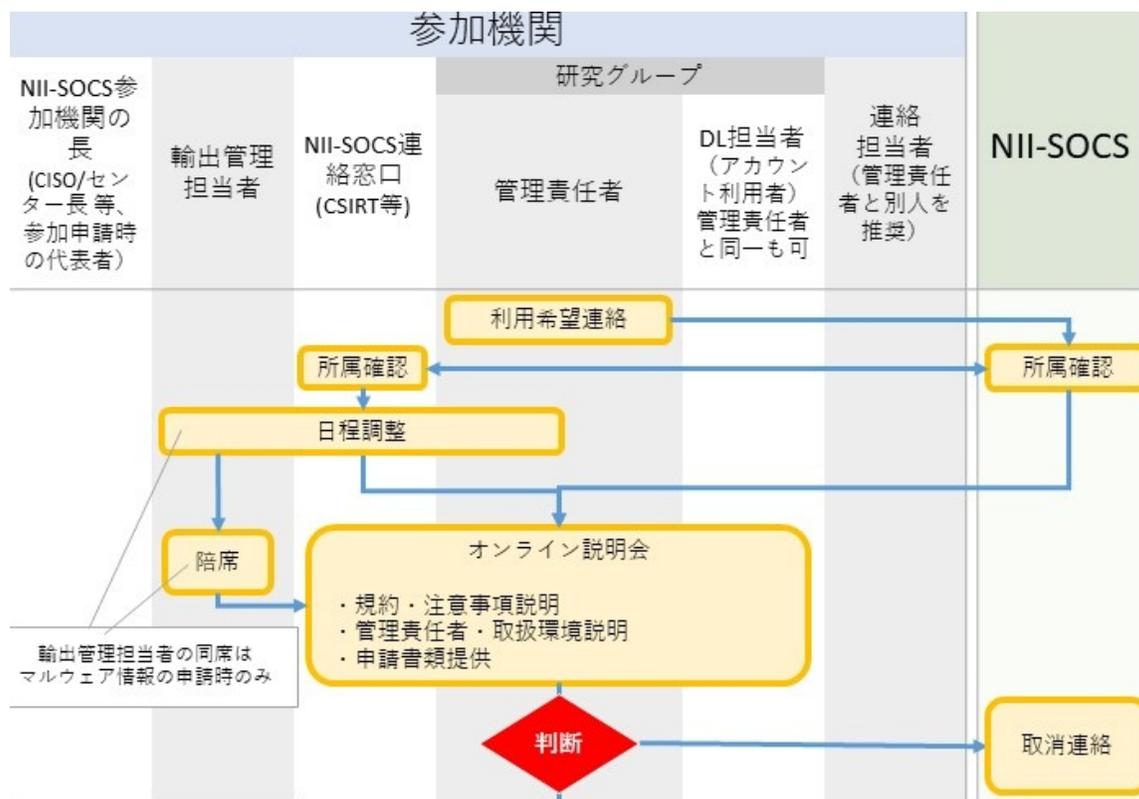
もっと利用者が増えて、研究成果がSINET加盟機関のセキュリティ向上につながると嬉しいということで、利用経験者から利用方法を説明

- 契約までの流れ
- 契約に必要な運用体制の準備
 - 嶋田研ではこうやっていますという事例紹介
- 研究実施から発表
- 個人的に特に気をつけている点
- さわってみた所感(研究事例)

契約までの流れ(2/3)

2. オンライン説明会(1時間程度)の開催

- マルウェア情報の申請も考えていて、まとめて説明してもらう場合は、説明会も輸出管理担当者にも同席してもらう必要あり
 - ...が、マルウェア情報の要求管理体制は厳しいので、「ついでに」な気持ちなら聞かない方が良い

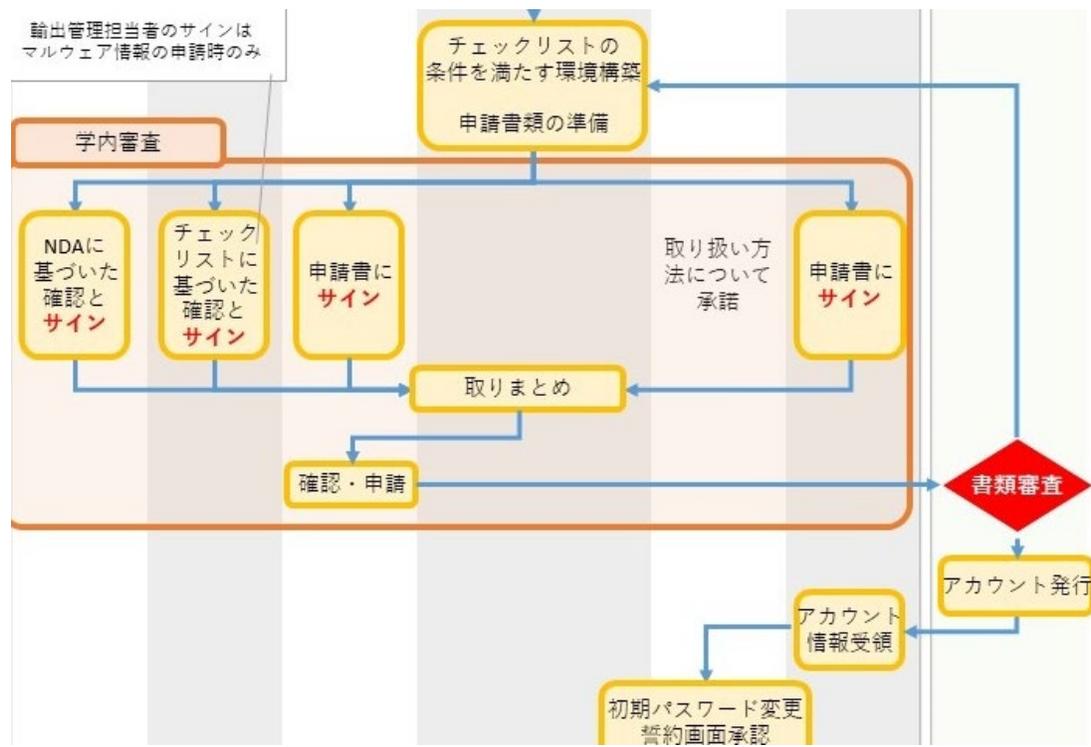


オンライン説明会で聞いたこと

- 説明自体は、申請手順/書類に沿った物
- 自分から「この形でOKか?」と確認する機会とすると良い
 - 専用ファイルサーバとか利用可能計算機とか認証とか
 - スパコンとか契約計算機での利用の確認
 - 内規作成上で迷っている所とか
- (マルウェア情報の必要管理体制は厳しくて、当分諦めることとした)
 - 確実な隔離ネットワークを要求されること(物理的にも)

契約までの流れ(3/3)

- 条件を満たす環境構築(内規、管理体制、など)
- 申請書類の提出と審査(出戻りあり)
- 審査が通ったらアカウント受領
 - データセット提供Webサイトから必要分をダウンロード



嶋田研で作成した内規(1/2)

以下の大項目に分けた上、各小項目を記した

- 冒頭の概要
 - 「研究室メンバであっても悪性通信研究班以外の人へのデータセットへのアクセスは固く禁止」を明記
- データセットとその危険性について
 - 「損害賠償も規定されている機微なデータセット」である点について言及
- データセットの管理
 - 研究責任者以外の連絡担当者の情報
 - 責任者不在時に連絡担当者経由でNII-SOCSから連絡の可能性の明記
 - 研究グループとしての定義
 - 流出などの事故や外部からの指摘に対しての連絡体制
 - 研究成果公表前のチェック(研究グループの2名以上で確認)

嶋田研で作成した内規(2/2)

● データセットの利用

- 利用開始時(初回レク、定期レク)
- 計算機上での扱い(扱いを許可された計算機、計算機上のアクセス権限設定、データの移動方法、利用記録を残し方)
- 研究成果の発表(内容を口外しない、成果公表前の管理者+1人以上による確認、発表に関連した指摘への対応)
- 利用終了時(データ消去記録と管理者への連絡、管理者による確認とアクセス権削除、作成した中間データの管理者への引き渡し)
 - 管理者は「論文の根拠となるデータの保持期間」の間、中間データをデータセット用ファイルサーバで保持
- データ流出の可能性があった場合の対応(不正アクセス発見、研究用PCの紛失/盗難、誤って許可されていない計算機への送付)

研究室外の教職員との体制構築

- 体制構築に必要な教職員(申請書に署名)
 - 研究グループ管理責任者
 - (アカウント担当者)
 - 連絡担当者
 - 嶋田研では常勤教職員が私だけなので、事務の総務係長に依頼
 - 各大学のNII-SOCS担当者
 - (マルウェア情報の場合は輸出管理担当)
- 研究室外から識別しやすい研究グループ名の設定
 - 例:「名古屋大学情報基盤センター嶋田研悪性通信研究班」
 - 研究室内でも「関係者以外は触っちゃだめ」と分かりやすいように設定

嶋田研でのデータ実体の運用

- NII-SOCSベンチマークデータダウンロードサイトは管理責任者(アカウント担当者は追加可能)しかアクセスできない
 - 研究室LAN内にNII-SOCSベンチマークデータ専用ファイルサーバを準備
 - 学生には個別にID/パスワードを渡す
 - 学生はデータのダウンロードや管理状況の記録をExcelテンプレートに残して、研究室ファイルサーバで管理
 - 専用ファイルサーバログ定期チェック&ちゃんと教員PCにもバックアップ(手動)
- 他のメンバにコンタクトがあることを考え、研究室内Wikiにも必要な情報は掲載
 - 現在/過去の利用者リスト、関連する教職員、現在の内規

学生用データ管理Excel

	DL日時	DLデータ範囲	DL先端末	嶋田DL確認	DL先端末消去日時
例	1970/1/1 0:00	2022/1/1-2022/1/31	嶋田PC(SAM: H11111111, シリアル: 123456)	1970/1/4 9:15	1970/3/20 12:00
1					

研究成果の発表時

- まず成果公表承諾依頼書を作成
 - タイトル、著者、公表先、概要などの論文情報
 - 「データの使い方」の記述、チェックリストに従ったチェックの実施
 - 全部でA4 1ページちょっと
- 著者1名の署名を入れた成果公表承諾依頼書を提出
 - 2週間前までに提出
 - 査読に出すのも公表扱い
 - 予稿の添付は不要(禁止)
 - 研究上の機密を先に知る形になってしまうので、送って来られると困る
- NII-SOCSで審議した上で承諾/追加確認等が来る

最終的なデータの使い方に責任を持つのは各研究グループ

- NII-SOCS側で研究上の機密に抵触する形で事前確認しないので

利用開始後の管理

- 毎年3月に継続利用の再申請の実施
- 毎年10月に管理体制のチェックリストによるチェック

参考: 嶋田の場合は利用希望連絡から実際に利用開始まで3ヶ月ほど

- 途中でお盆をはさんだり、学内側で必要となる担当者や署名の依頼方法を利用希望連絡後から確認したので、かなりゆっくり
- ちゃんと事前準備をした上で進めれば、1ヶ月ちょっと程度で利用開始できるのでは?

規約上で個人的に気をつけている点

- 一般的な注意事項の他に、個々のデータからの固有情報の特定や公表(偶然に見つかってしまった)の禁止

第6条 利用者への禁止事項と対応措置については、次の各項のとおりとします。

1. 利用者は、本約款で許諾される場合を除き、次の各号の一に相当する、又はそれと同等の行為を行ってはなりません。

---snip---

七. コンテンツから、送信元・送信先 IP アドレス等、通信機器やその使用者の固有情報を特定することを目的とした行為

八. 前号の目的の有無にかかわらず、コンテンツから特定された通信機器やその使用者のIPアドレス等の固有情報の公表

九. 国立情報学研究所の承諾を得ない運用連携サービスの機能や検知手法に関する情報の公表

十. 国立情報学研究所及び機器製造元の承諾を得ない運用連携サービスで使用する機器の機能や検知手法に関する情報の公表

十一. 本サイトについて、リバースエンジニアリング, 逆アセンブル, 逆コンパイル又はソースコードの抽出を行うこと

少しNII-SOCSベンチマークデータを触ってみた印象

- 1日あたりのデータがかなり多い印象
 - 1日分のデータ(展開済み)の中央値は5GB前後になる印象
 - Kyoto 2016 Datasetは1日分(展開済み)で数十MB程度
 - 悪性/良性(無害)のラベルの付与された通信データセットとしてだけでなく、単純なビッグデータとしての扱いも面白そう
- プローブ系の通信も大量に入った、本当に生の通信データセットという印象
- 良くある悪性/良性(無害)のラベルのついた通信データセットと比較して、通信のバリエーションが多い印象
 - Kyoto 2016 Dataset(ハニーポットで取得)と、とりあえず悪性/良性(無害) 識別器を作って比較(後スライド)すると識別しづらい

提供されるNII-SOCSベンチマークデータの仕様

- SINETと商用系の間通信を観測するNII-SOCSが取得するデータの一部を匿名化して提供
 - /24(IPv4)と/64(IPv6)のブロックをランダムに10ブロック程度ずつ選択、1週間単位で切り替え
 - 想定以上の流量となった場合は、セッションの一部が欠落する可能性も
 - そもそも、通信データセットでセッションの完全性を保証している事例はまず無い
- Kyoto 2016 Dataset準拠+αで匿名化
- データセット提供Webサイトで1日ごとのtarballで提供される
 - サマライズ+匿名化された通信セッション情報、検出サマリ、事後検証結果、のファイルが含まれる[1]
 - 1日分のデータtarballのサイズは500MB-2GBあたり、1GB前後が中央値という印象
 - 現状で1年半分が提供されている

[1] 「大学間連携に基づく情報セキュリティ体制の基盤構築」事業における研究用データの提供について / 研究用データの解説
<https://meatwiki.nii.ac.jp/confluence/pages/viewpage.action?pageId=59022903>

通信セッションデータの匿名化方法

- かなりがっつりと匿名化されている

ランダム化処理方法

1. 毎日、以下の条件でトラフィックデータを抽出(pcapファイルの生成)
 - a. 参加機関のIPアドレス領域から/24(IPv4)または/64(IPv6)のブロックをランダムに複数選択。(IPv4、IPv6ともに10ブロック程度、このセットを7日間使用する。)
 - b. 00時00分00秒から23時00分00秒の間からランダムに30分間の枠を二つ選択。
2. 観測対象時刻のトラフィックデータのタイムスタンプを当該日の0時0分0秒から0時29分59秒と12時0分0秒から12時29分59秒に振り直す。
3. 送信元IPアドレス/受信先IPアドレスをIPv6形式のランダムなIPアドレス領域に振り直す。
 - a. IPv4は/24、IPv6は/32の範囲内のIPアドレスの連続性を維持する。
 - b. ランダム化処理に使用するseedは定期的に変更する。
4. ポート番号についてはwell-known port(1024未満)はそのままとし、それ以外をランダムな値に振り直す。
 - a. ランダム化処理に使用するseedは7日おきに変更する。
5. 他はKyotoData2016[1][2]に準拠した統計データとし、ペイロードやDNS名は含まない。

NII-SOCSベンチマークデータの検知ラベルの仕様

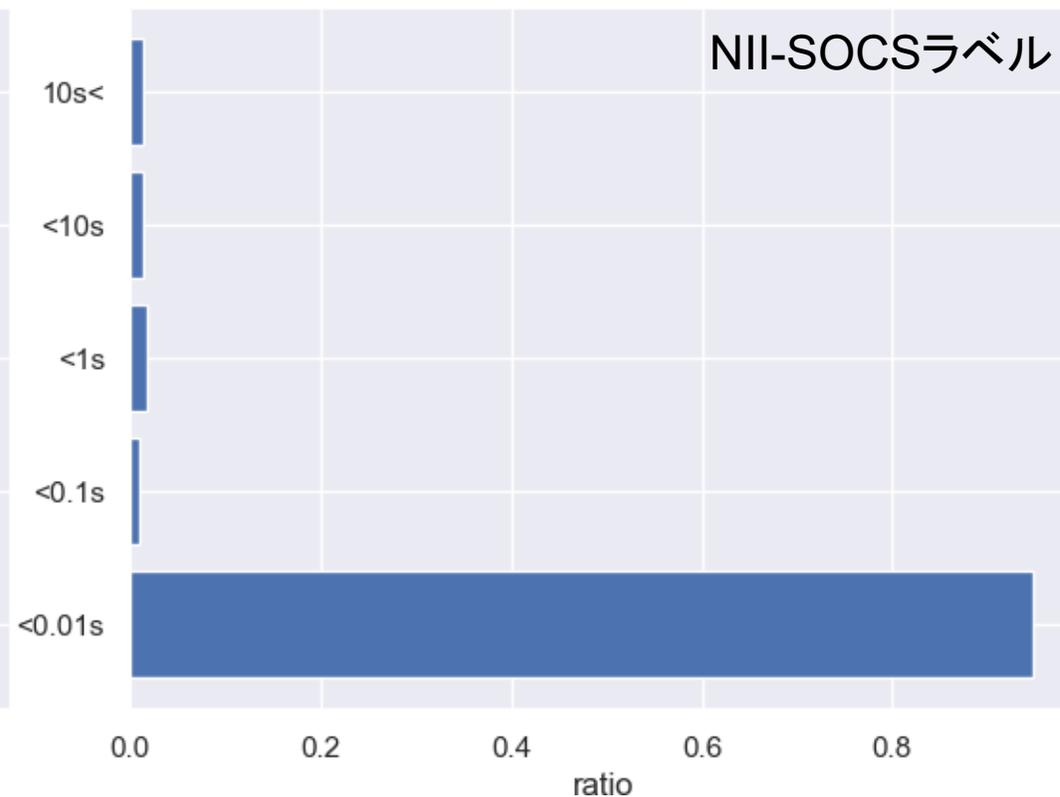
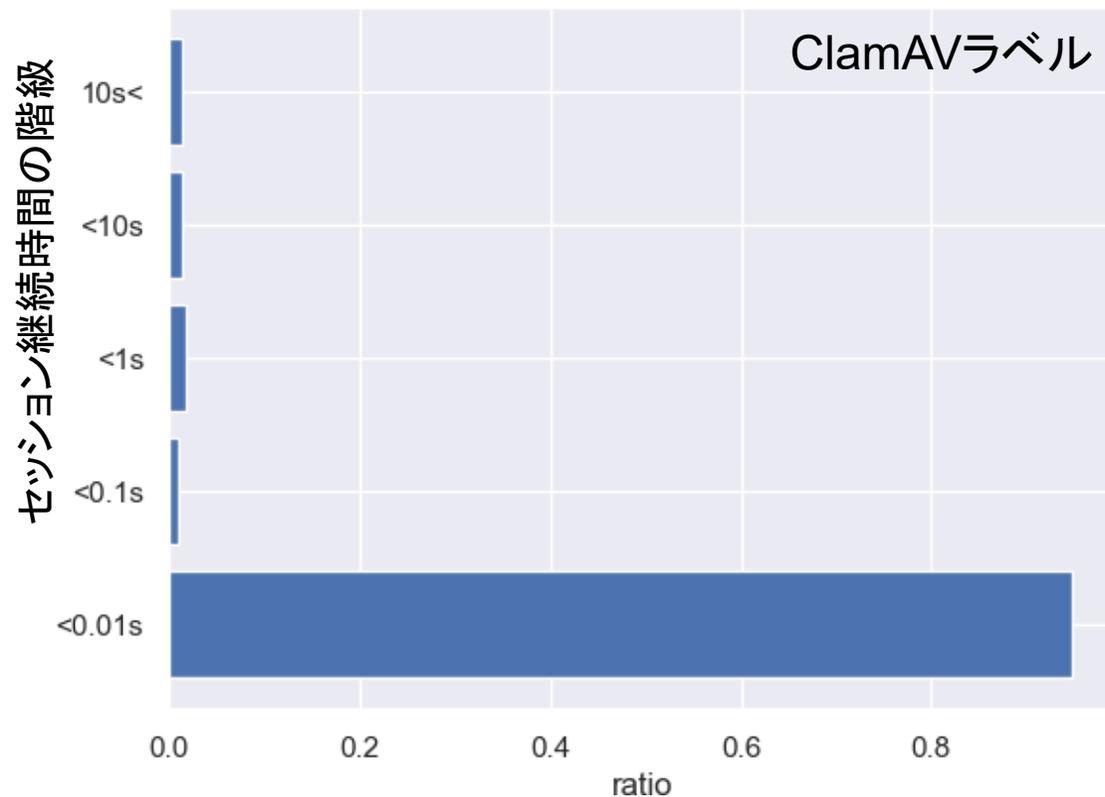
- Snort / ClamAV / NII-SOCSの検知ラベル
 - NII-SOCSは複数の検知システムの検知
 - 当初は反応が無いけど、新しい検知ルールが出た後に、再度検知をかけたら反応があったら「ZERODAY」に入れられる
 - SnortとClamAVについては採取日翌日に検査、以後7日おきに8回(採取57日後)の検査
 - Snortの場合は採取日翌日は未検知、2-5回目(8-36日目)検知、6-8回目(43-57日目)検知のラベル付きがZERODAYへ
 - 無償版Snortのため、30日程度シグネチャ提供が遅れることがある
 - ClamAVの場合は採取日翌日は未検知、2回目以降検知のラベル付きでZERODAYへ
- 詳細はNII-SOCSのMeatWiki(open)参照[1]

[1] 「大学間連携に基づく情報セキュリティ体制の基盤構築」事業における研究用データの提供について / 研究用データの解説
<https://meatwiki.nii.ac.jp/confluence/pages/viewpage.action?pageId=59022903>

良性(無害)通信のセッション継続時間

- ClamAVラベル、NII-SOCSラベルともども、良性(無害)のラベルがついている物の95%程度が0.01秒未満

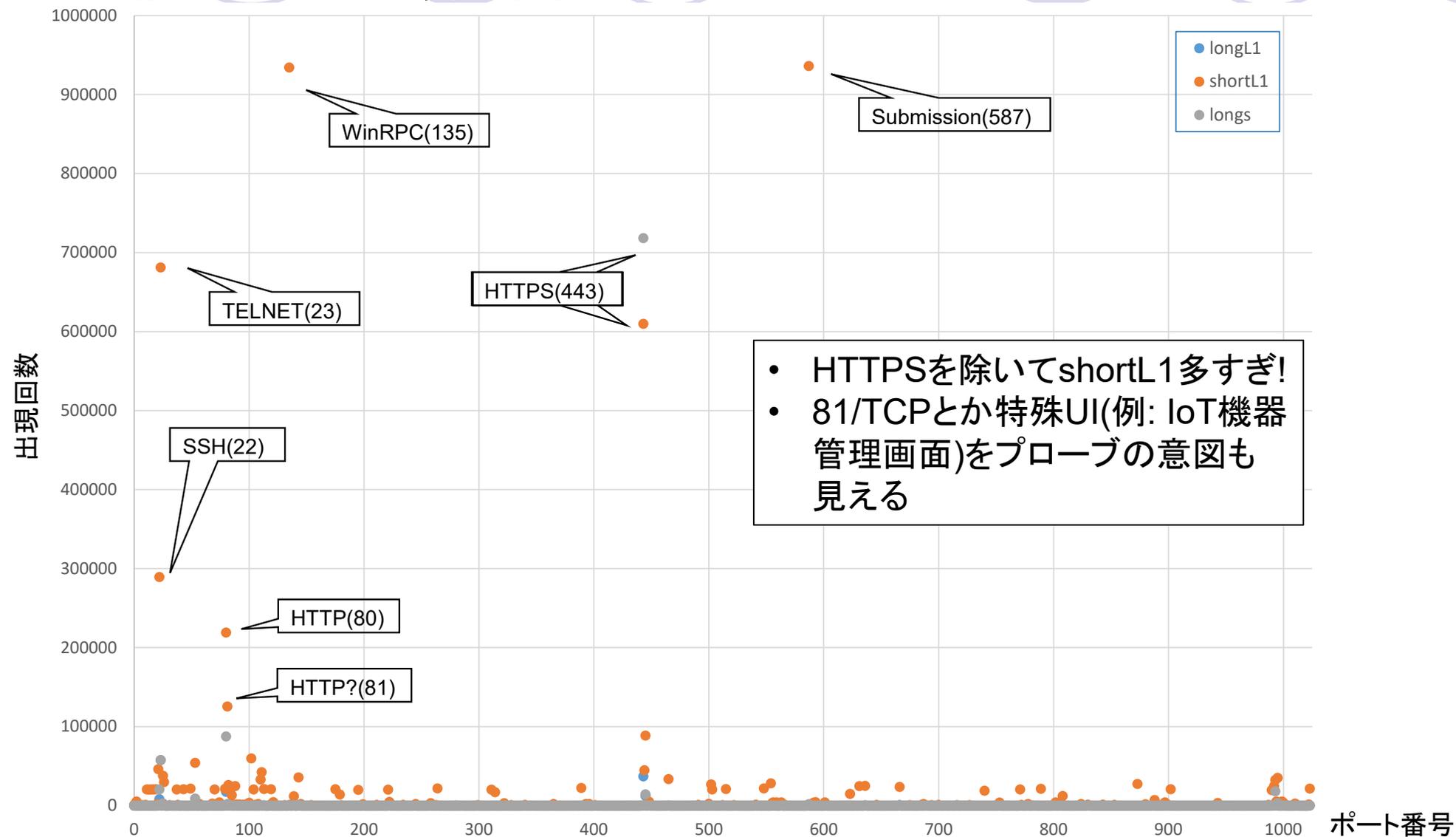
→プローブ系の通信が大多数を占めていそう



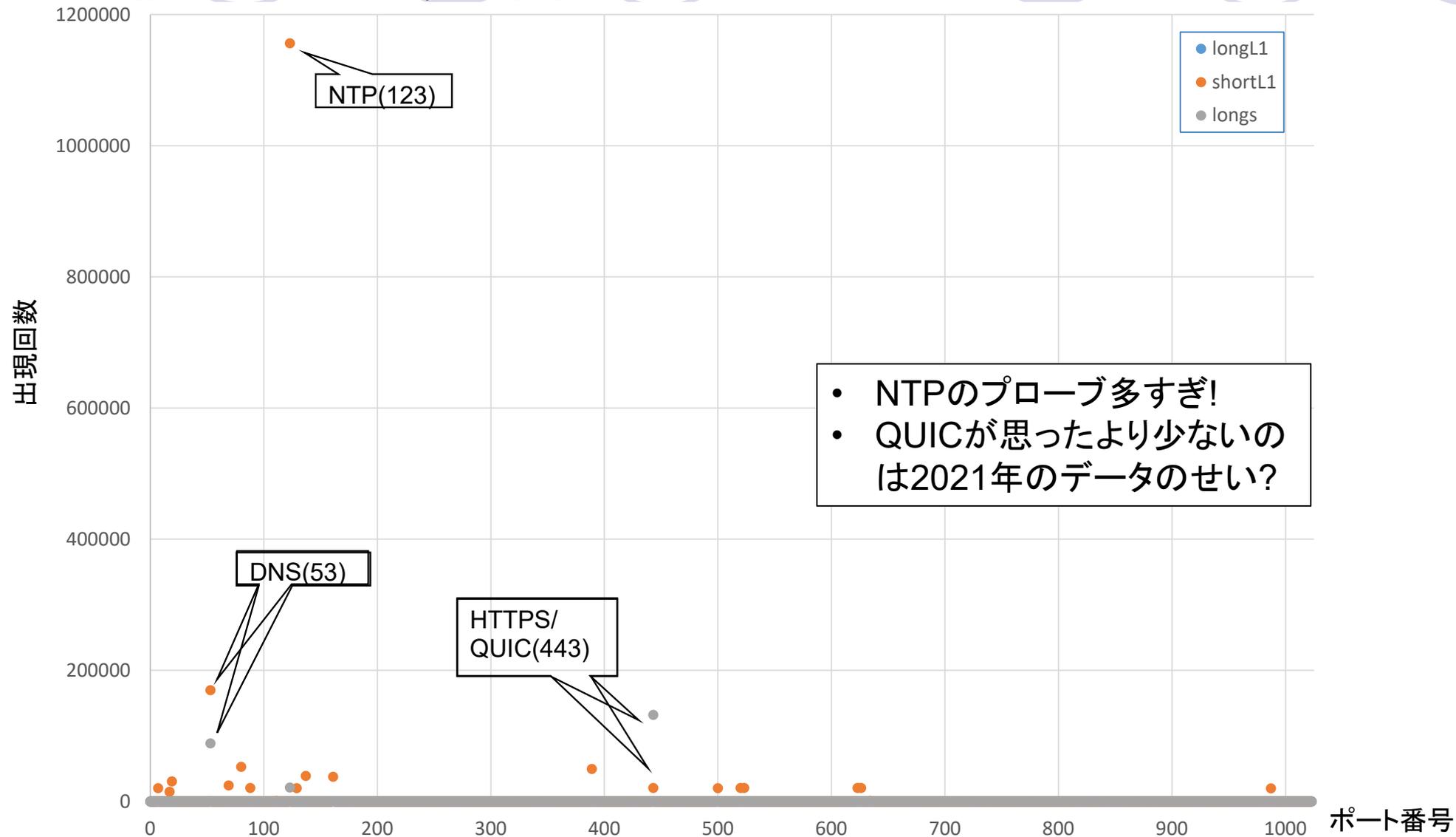
プローブ系の通信(short L1)の除外の試み

- 現状、「KyotoData 2016基準で、セッション終了時の接続状態がL1に属し、なおかつ、通信の継続時間が0.01秒以下(short L1)」をプローブ系通信として除外した方が良さそう
- 「セッション終了時の接続状態L1」は以下の2種類
 - 接続状態S0: 接続試行が見られない
 - 接続状態REJ: 接続が拒否された
- 参考(他の接続状態)
 - L2: 接続が正常に終了した、または接続中
 - L3: 接続開始側に問題が発生
 - L4: 接続応答側に問題が発生
 - L5: 接続応答側からSYN-ACKパッケージが開始側に届かない
 - L6: 接続開始側からのSYNパッケージが応答側に届かない

ある日のshort L1/long L1/longs(0.01秒以上全て)のTCP通信の出現頻度



ある日のshort L1/long L1/longs(0.01秒以上全て)のUDP通信の出現頻度



short L1除外も含めた良性/悪性通信の識別器作成と評価

- とりあえず、特徴量(20種類)を用いてLightGBMで識別器を作ってみた
 - 悪性/良性(無害)ラベルのついた1万件ずつ抽出
 - 半分(1万件)で学習、残りの半分(1万件)を識別
 - 結果([1]に一部掲載)
 - 前処理しなくてもKyoto 2016 Dataset(ハニーポット)よりも識別が難しい
 - short L1除外の前処理を入れるとさらに評価値が落ちる
- 悪性通信検知の研究用データセットとしてやりがいがある!

	Kyoto2016 Dataset(参考)	short L1 除外無し	short L1 除外あり
正解率	0.9867	0.9490	0.7964
適合率	0.9877	0.9411	0.7594
再現率	0.9856	0.9580	0.8678
F1 score	0.9867	0.9495	0.8100

まとめ

情報基盤センター所属の研究者の嶋田として

- NII-SOCSベンチマークデータを用いたセキュリティ研究が盛り上がると嬉しい
- 研究成果がSINET加盟機関のセキュリティ向上につながると嬉しい

なお

- 利用規約違反には気をつけましょう
- 研究用としてかなりがっつり匿名化されている点に注意
 - IPアドレスは/24(IPv4)、/32(IPv6)の範囲のみ連続性を保って匿名化
 - 同じIPアドレスブロックからのデータ取得は7日間
 - 観測日時も移動(連続性は担保)、high portはランダム化
→1ヶ月レベルの統計とかをベースに研究するのが良さそう
 - 日本の全SINET加入機関を対象とした攻撃の検知としても