

敵対的生成ネットワークを利用した 疑似トラヒックデータ生成手法

長岡技術科学大学 大学院工学研究科

電気電子情報工学専攻

山際 哲哉

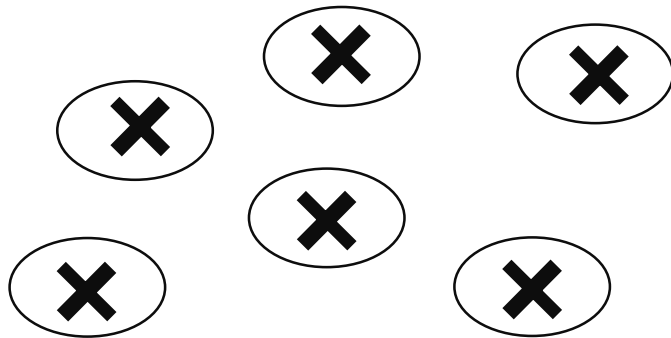
渡部 康平

中川 健治

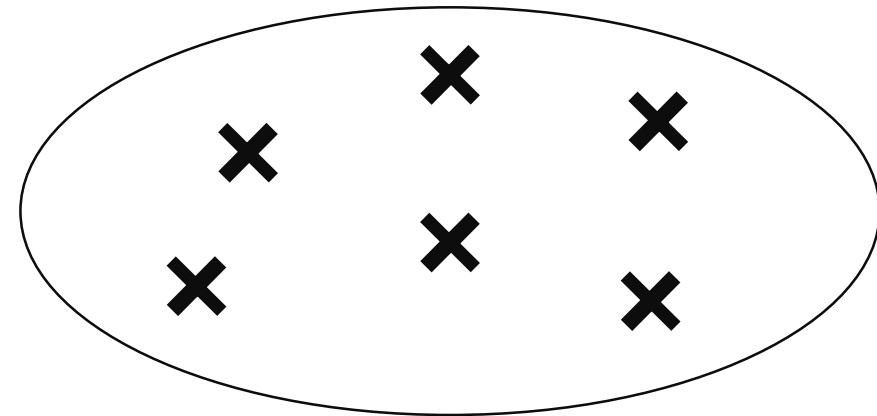
- 近年，ユーザーが利用する端末が多様化し，ネットワークトラヒックの性質も多様化している
- ネットワークを快適に利用できる最適なシステム提供が必要
- その際トラヒックジェネレータを使って，シミュレーションやテストを行う
- 例：キャパシティプランニング
 - ITシステムの構築の際に使われる
 - あるトラヒックを処理するのに，どのくらい増強すれば良いか
- 公開されているデータセットが少ない
- リアルなトラヒックを作るのは困難
 - 統計学的知識やパラメータ設定

- 機械学習のGANを使って、誰でも簡単に多くのリアルなトラヒックを生成する
- 既存のトラヒックジェネレータとGANで作るデータの違い

既存のトラヒックジェネレータで
作る似たデータ

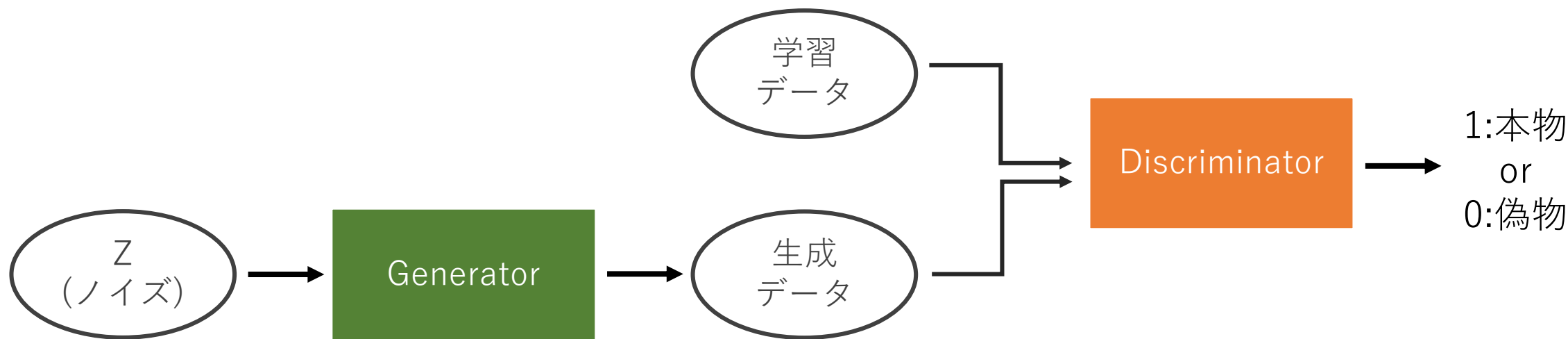


GANで作る似たデータ

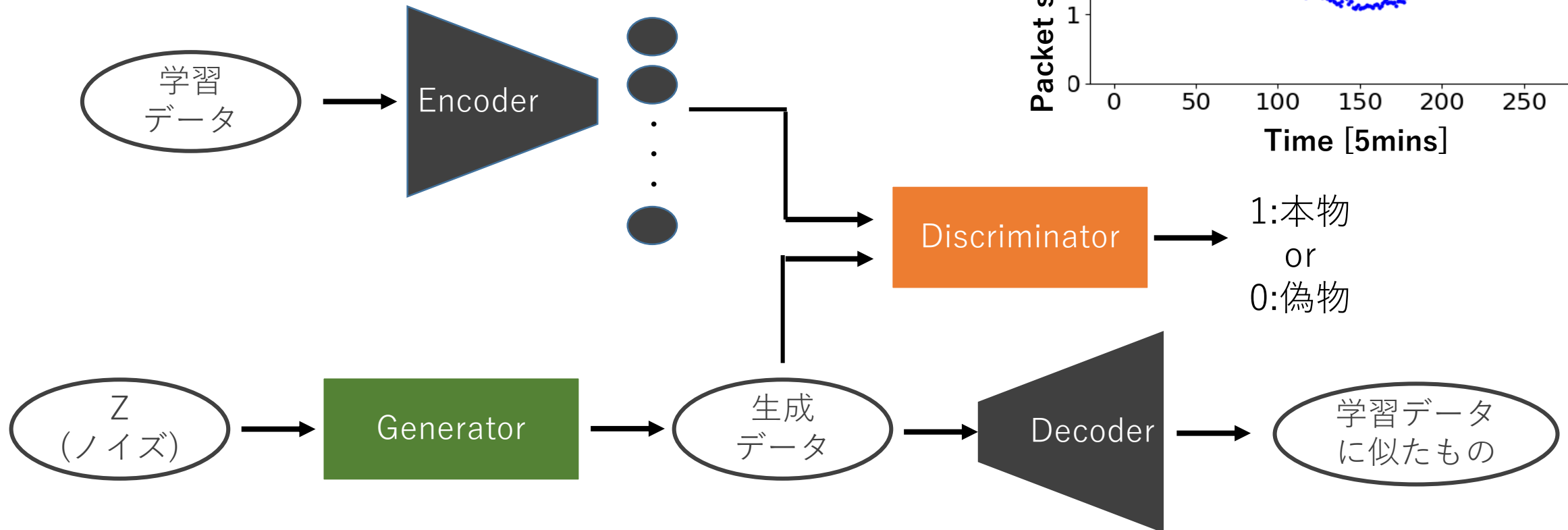


× : サンプル
○ : 生成する範囲

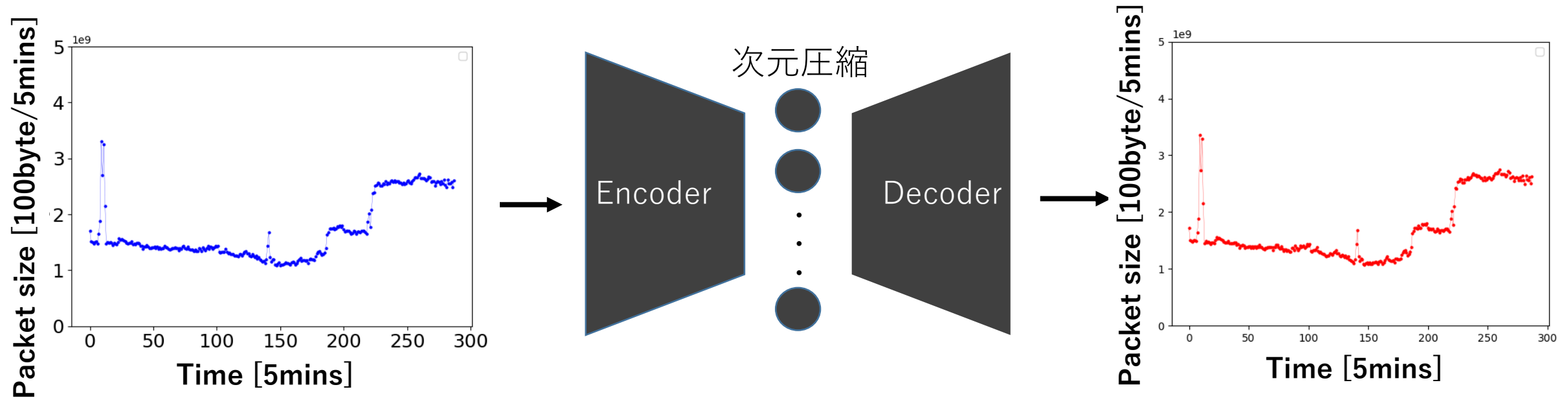
- 本物らしい画像生成の教師なし学習モデル
- 識別器(Discriminator)と生成器(Generator)の2つで構成されている
 - Discriminatorは学習データと生成データを判別する
 - 学習データなら1, 生成データなら0となるように学習する
 - GeneratorはDiscriminatorを騙せるほど類似したデータを生成する
 - Discriminatorの出力が1に近づくように学習する



- Encoderを用いてGANに学習
 - 学習データは時系列(各次元はトラフィック量)
 - データにバースト性があり離散値のため通常のGANの学習が困難

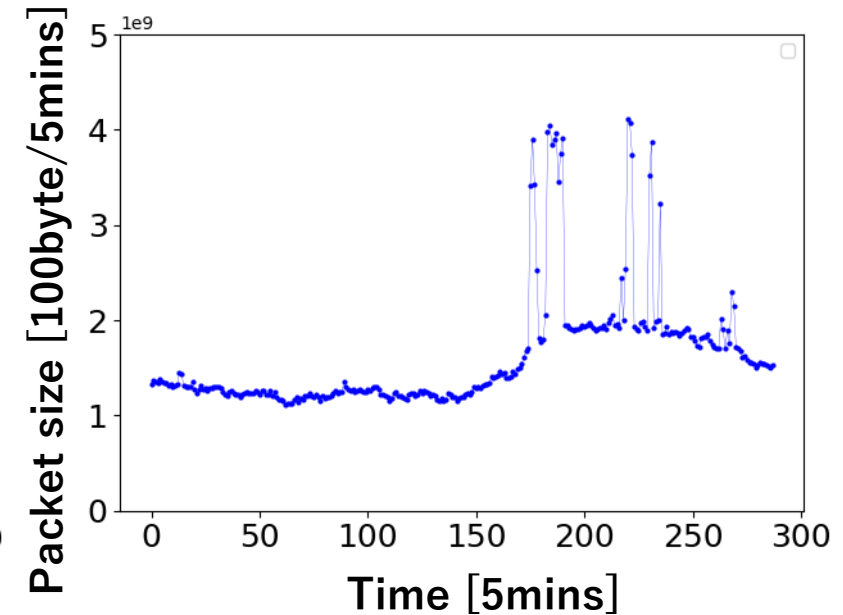
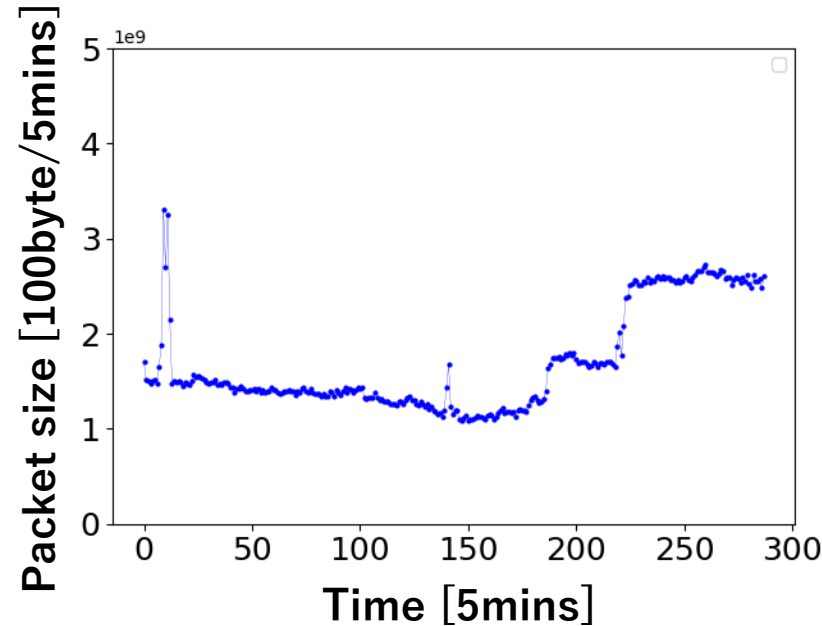
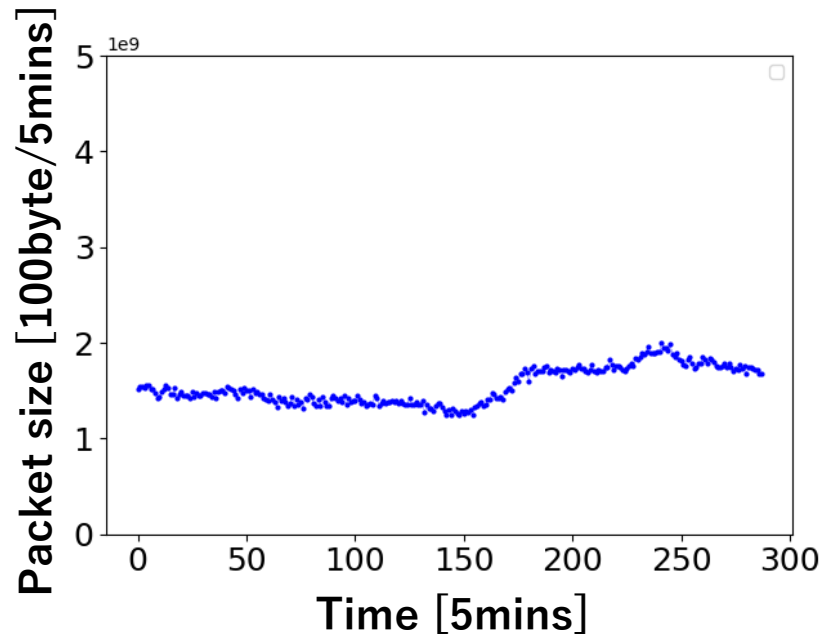


- データの次元圧縮として使われている
- Encoder, Decoder
 - データの次元圧縮, データの復元
- Encode後の次元は各学習データの特徴量を示している



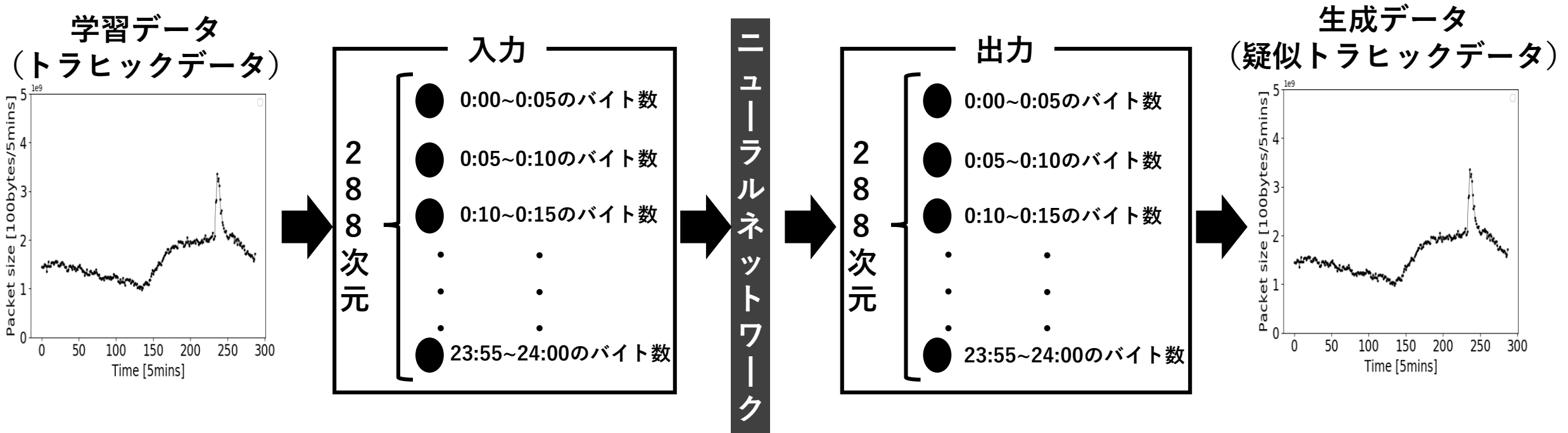
- 次元圧縮することでGANの生成範囲を狭めることができる

- 公開されているトラフィックデータ[2]を学習させる
 - 1日の5分間毎に何バイト通信されたかを表したデータ
 - 1日分を1つの学習データ
 - 168日分



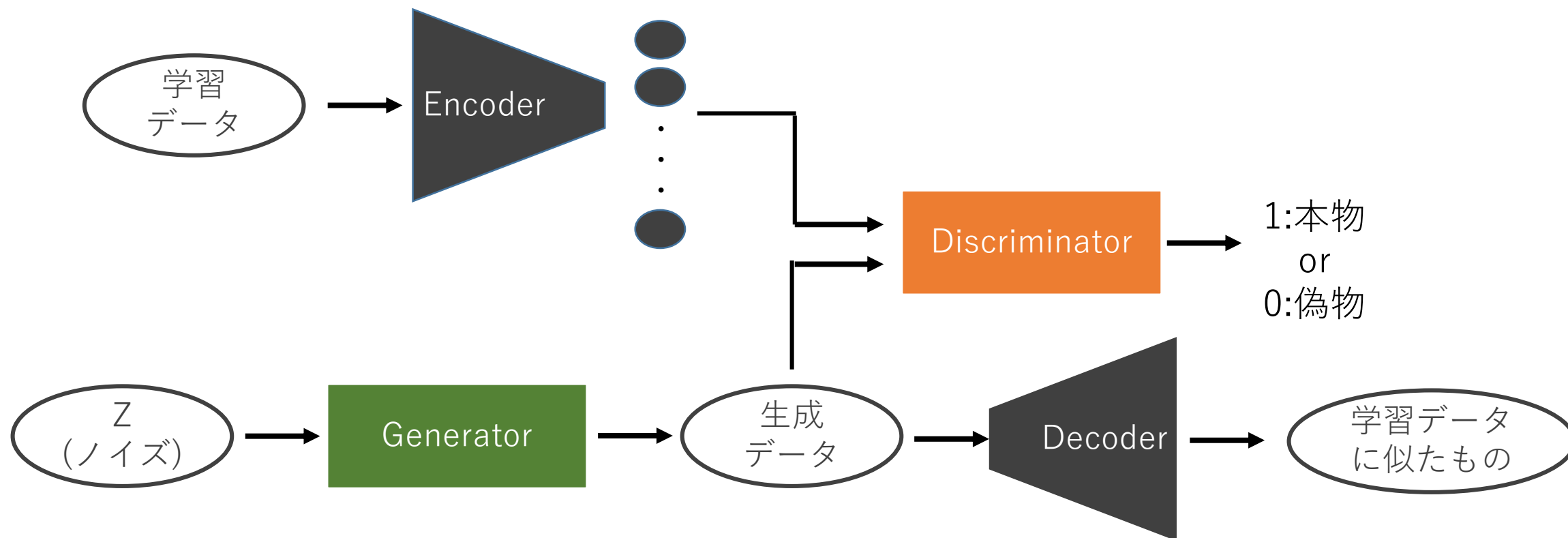
- トラフィックデータをGANによって生成する

- トラフィックデータをニューラルネットワークに入力し，疑似トラフィックデータを出力する
 - 入力：時間間隔毎のバイト数を成分に持つベクトル
 - 出力：疑似トラフィックデータを表す同様のベクトル

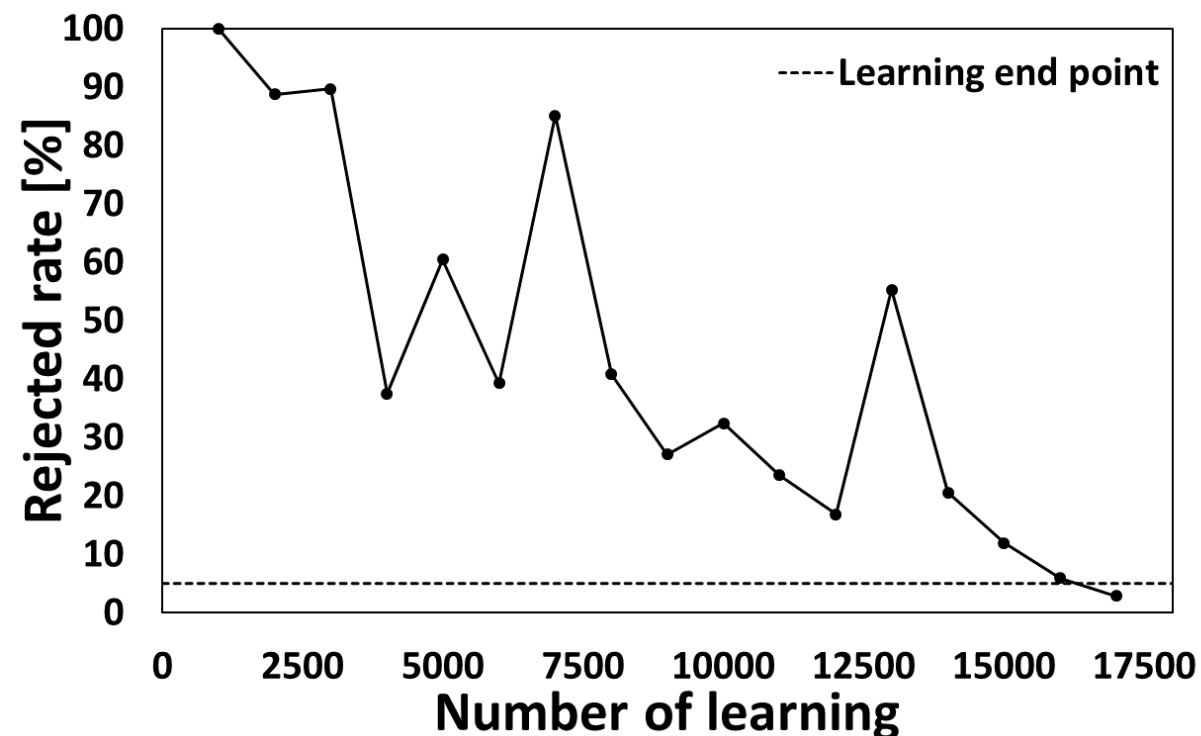


• 学習方法

1. トラフィックデータをEncoder, Decoderに学習させる
2. EncoderとGANを組み合わせて学習させる (GANのみを学習)
3. Generatorからの生成データをDecoderに入力し, 疑似トラフィックデータを生成する



- KS(コルモゴロフ-スミルノフ)検定
 - 仮説「標本XとYが同一の母集団の確率分布から発生」
 - 棄却された場合：確率分布が一致していない
- 5分間毎のバイト数における学習データと疑似トラフィックデータをそれぞれ168サンプルで検定する
 - 100セットを行う
- KS検定で棄却される割合が5%を下回ったとき，GANの学習を終了させる
 - 有意水準を5%に設定しているため

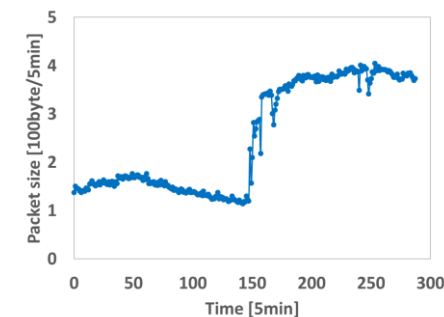
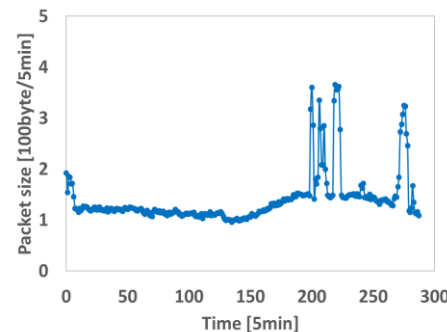


- 各ニューラルネットワークの構成

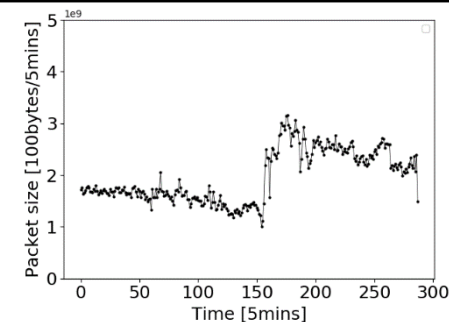
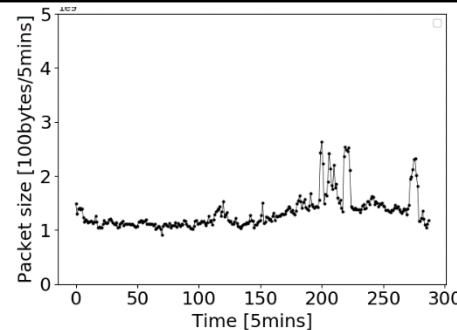
	Generator		Discriminator		Encoder		Decoder	
Layer	Units	Act.	Units	Act.	Units	Act.	Units	Act.
Input	2	-	5	-	288		5	-
Hidden	10	LReLU	50	LReLU	400	LReLU	50	LReLU
Hidden	30	LReLU	40	LReLU	200	LReLU	100	LReLU
Hidden	40	LReLU	30	LReLU	100	LReLU	200	LReLU
Hidden	50	LReLU	10	LReLU	50	LReLU	400	LReLU
Output	5	-	1	Sigmoid	5	Tanh	288	-

- 生成結果

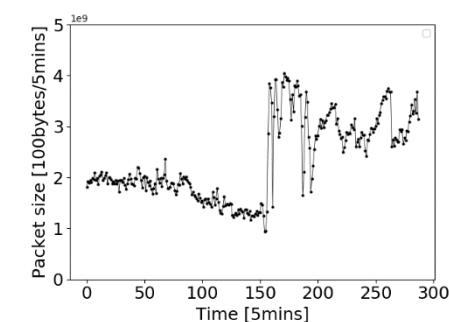
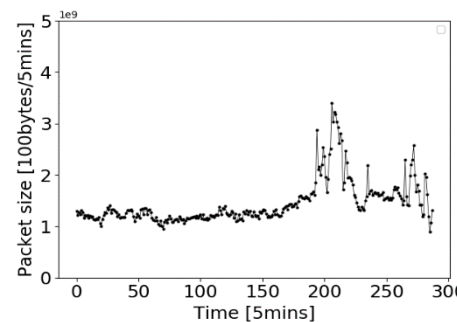
学習データ



生成データ
(類似)



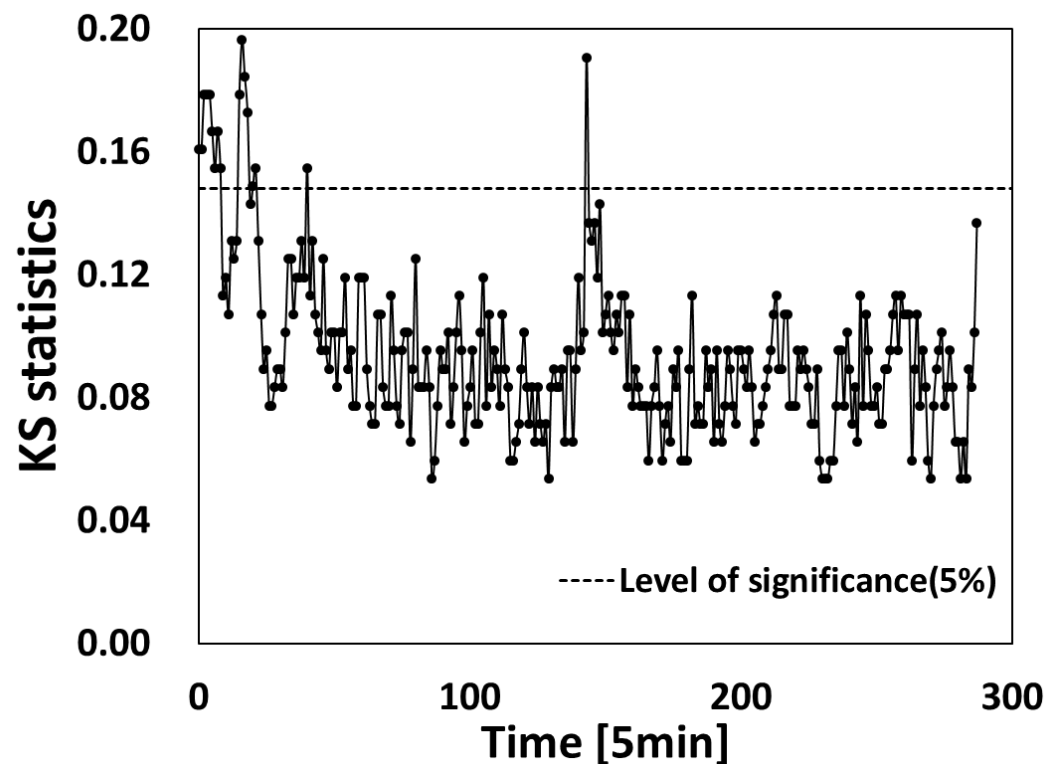
生成データ
(一部類似)



- 学習データに類似しているデータを生成していた

- KS検定

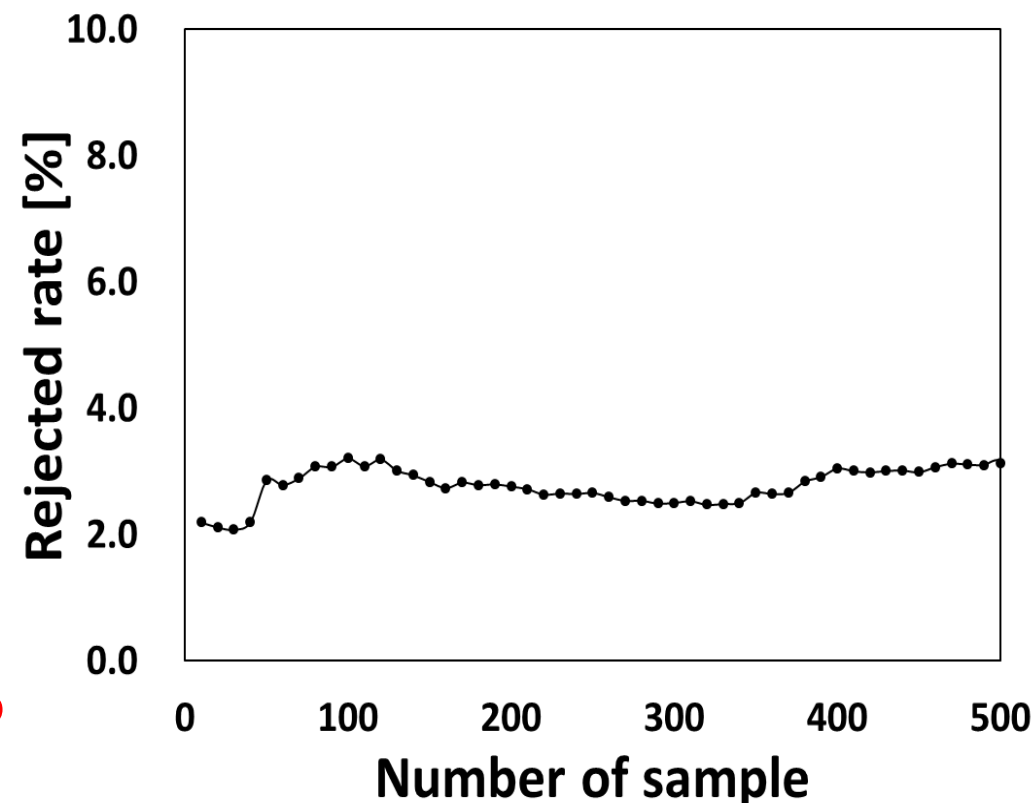
- 仮説「標本XとYが同一の母集団の確率分布から発生」
- 棄却された場合：確率分布が一致していない



- 生成データは本物のデータの分布とかけ離れていない

- GANはランダムな生成のため，必ずしもKS検定で棄却されない訳ではない
- GANで500セット生成し，棄却される割合を確認したところ，3.1%に収束した
- 有意水準5%に設定しており，実データで検定を行った場合，95%棄却されないが，GANの場合96.9%棄却されていない

- GANの生成は本物のデータとほぼ同等
→ 本物のデータの再現がほぼできている



- まとめ
 - Encoderを用いることでGANを使って、トラフィックデータ（学習データ）に類似したデータを作り出せることを確認した
- 今後
 - より精度の高いGANの構築
 - 多様性のあるトラフィックデータの生成
 - 時系列での評価