



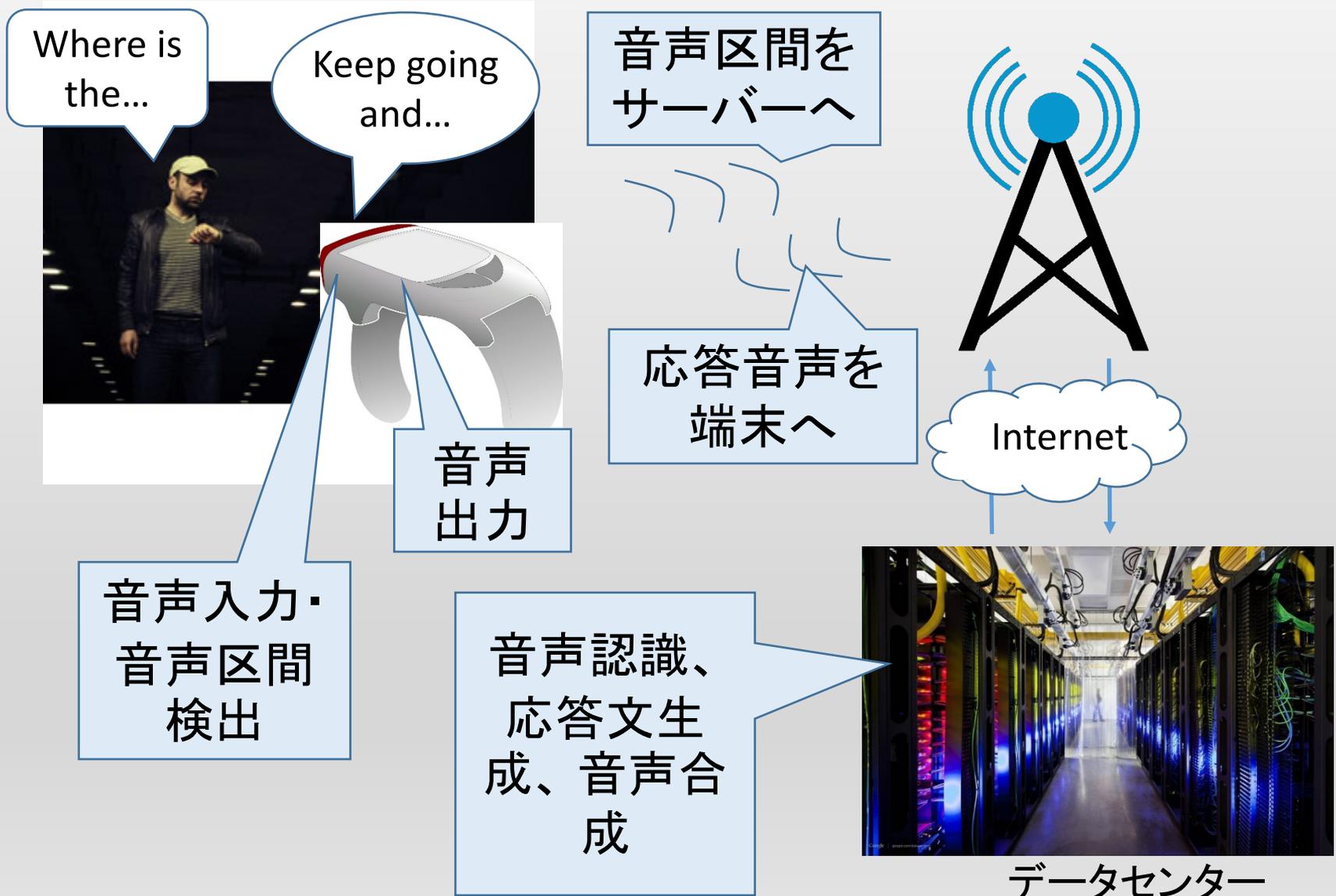
大規模進化計算による 音声認識システム開発の自動化

篠崎隆宏

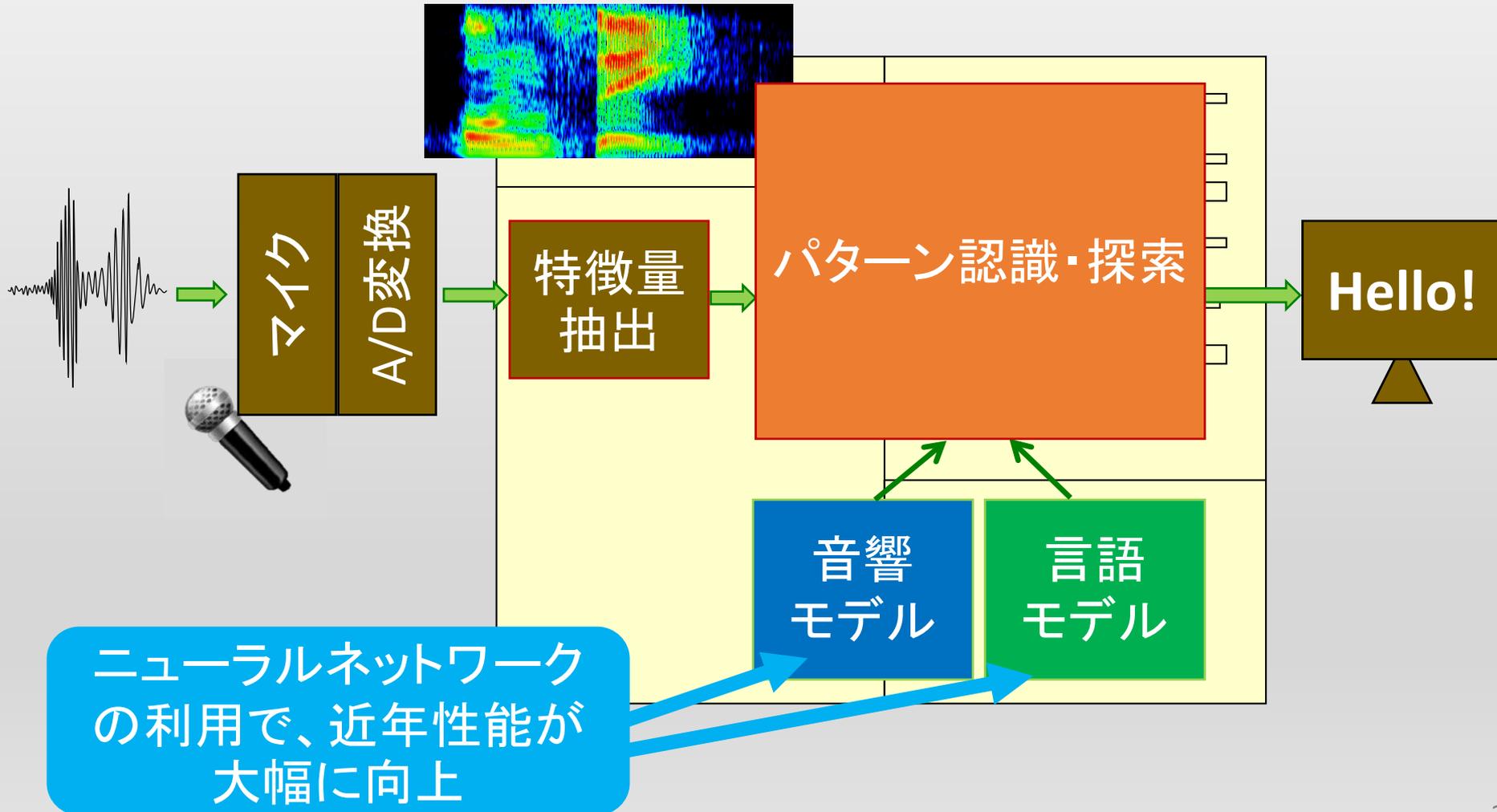
東京工業大学 情報通信系

www.ts.ip.titech.ac.jp

音声認識の普及

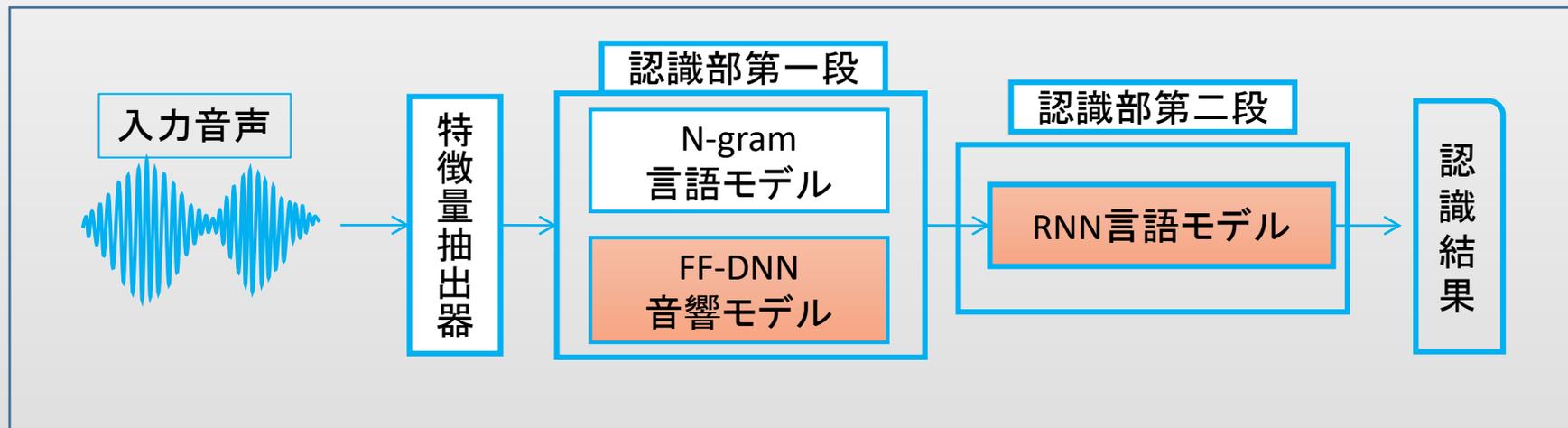


音声認識の仕組み

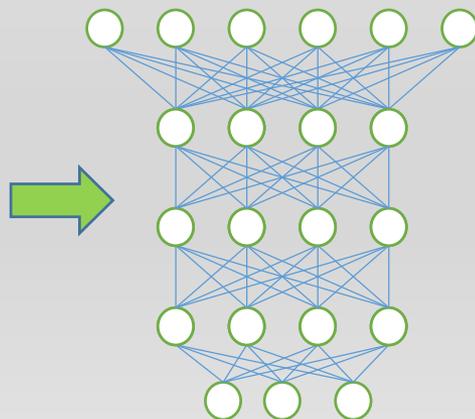


音声認識システムの構築

① システム設計



② 音響・言語モデルの学習



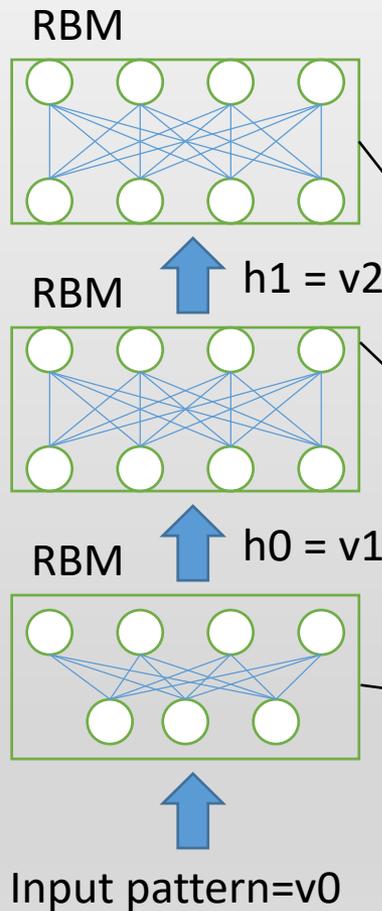
③ 認識性能の評価

単語誤り率(WER)

$$= \frac{\text{置換誤り} + \text{挿入誤り} + \text{削除誤り}}{\text{正解単語数}}$$

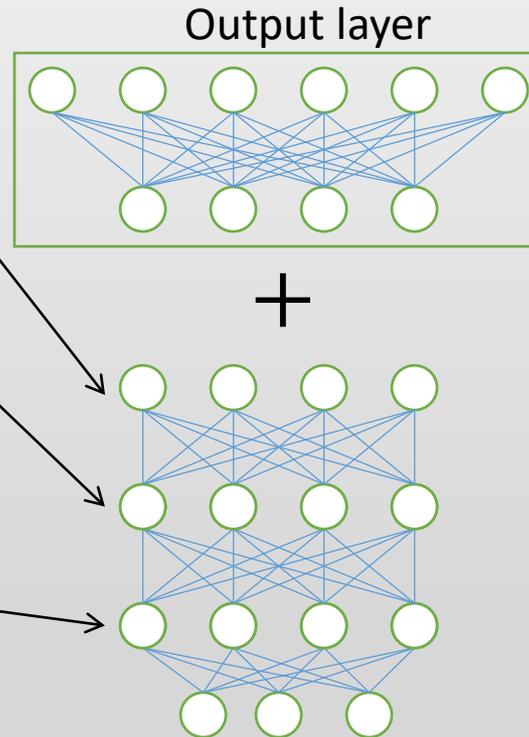
DNN音響モデルの学習プロセス

① RBMの教師なし学習

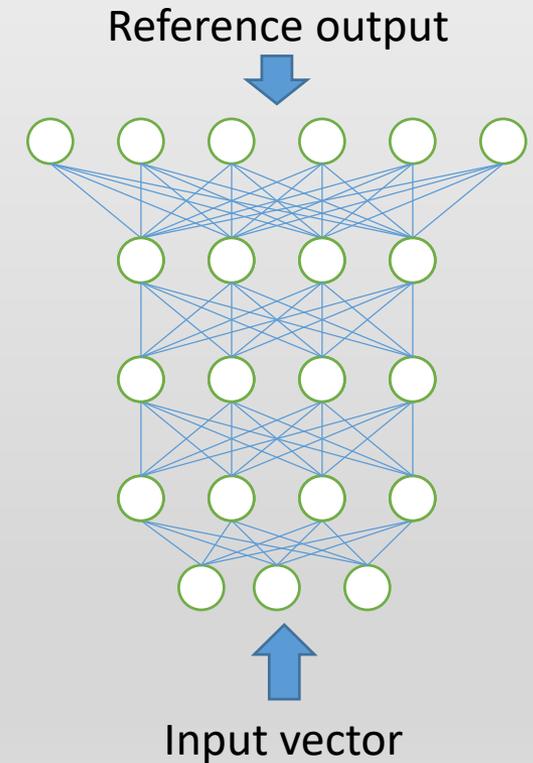


プレトレーニング

② DNNの初期化



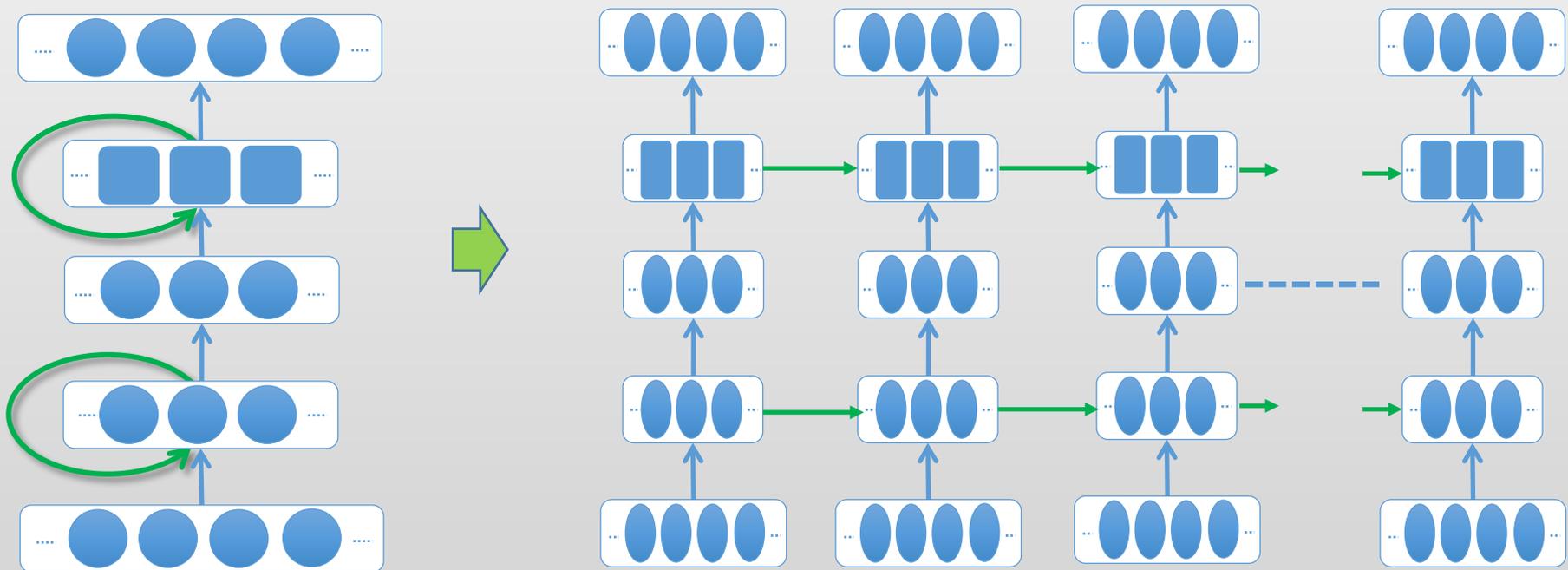
③ バックプロパゲーションによる教師付き学習



ファインチューニング

RNN言語モデルの学習プロセス

① RNNのループを時間方向に展開

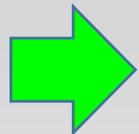
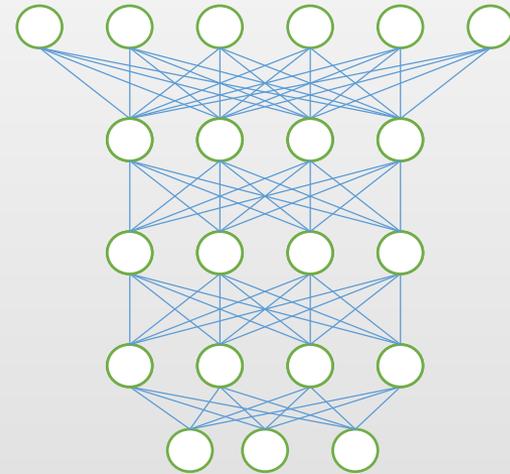


② 時間展開バックプロパゲーション (BPTT)により教師付き学習

音響・言語モデルの学習の計算規模



大規模な音声データ
(数百～数千時間)



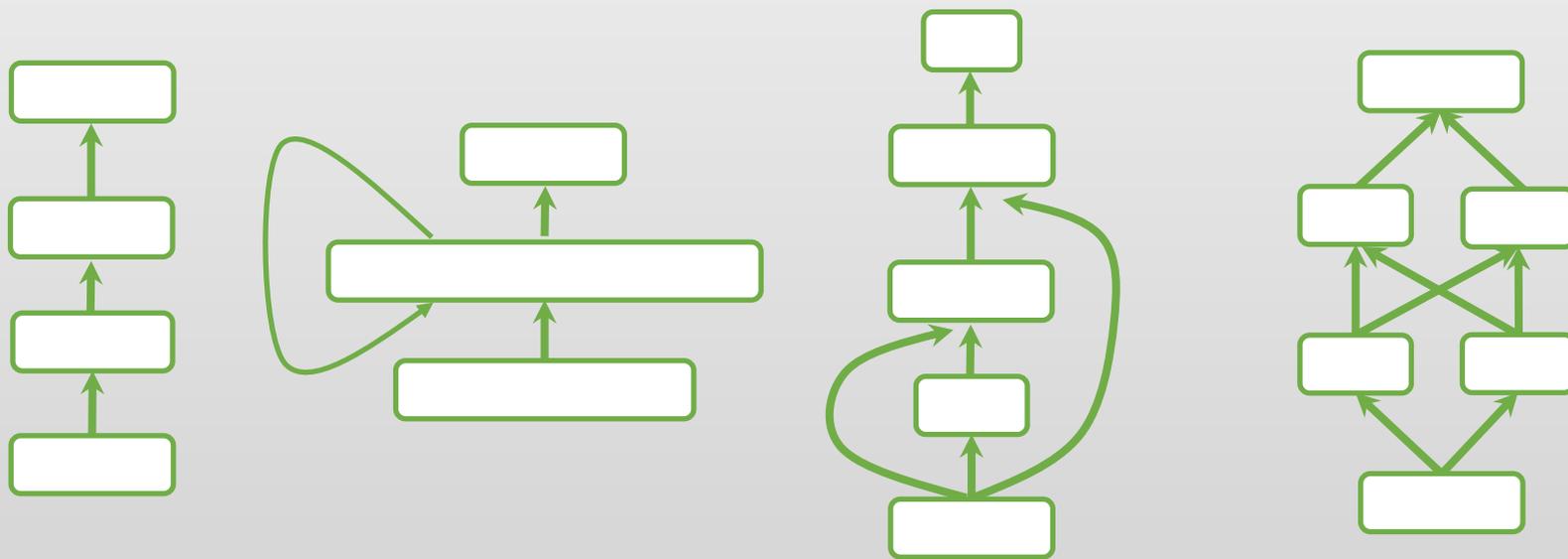
最新GPGPUを用いて数日～数か月



音声認識システム構築の問題点

相互依存する多数の調整要素(メタパラメタ)が存在

- ・ニューラルネットワークの構造
(隠れ層数、各層のユニット数、ユニット種別、接続構造等)



- ・学習条件
(学習率、モーメントム、ミニバッチサイズ等)

ノウハウと試行錯誤の世界

従来アプローチと提案アプローチ

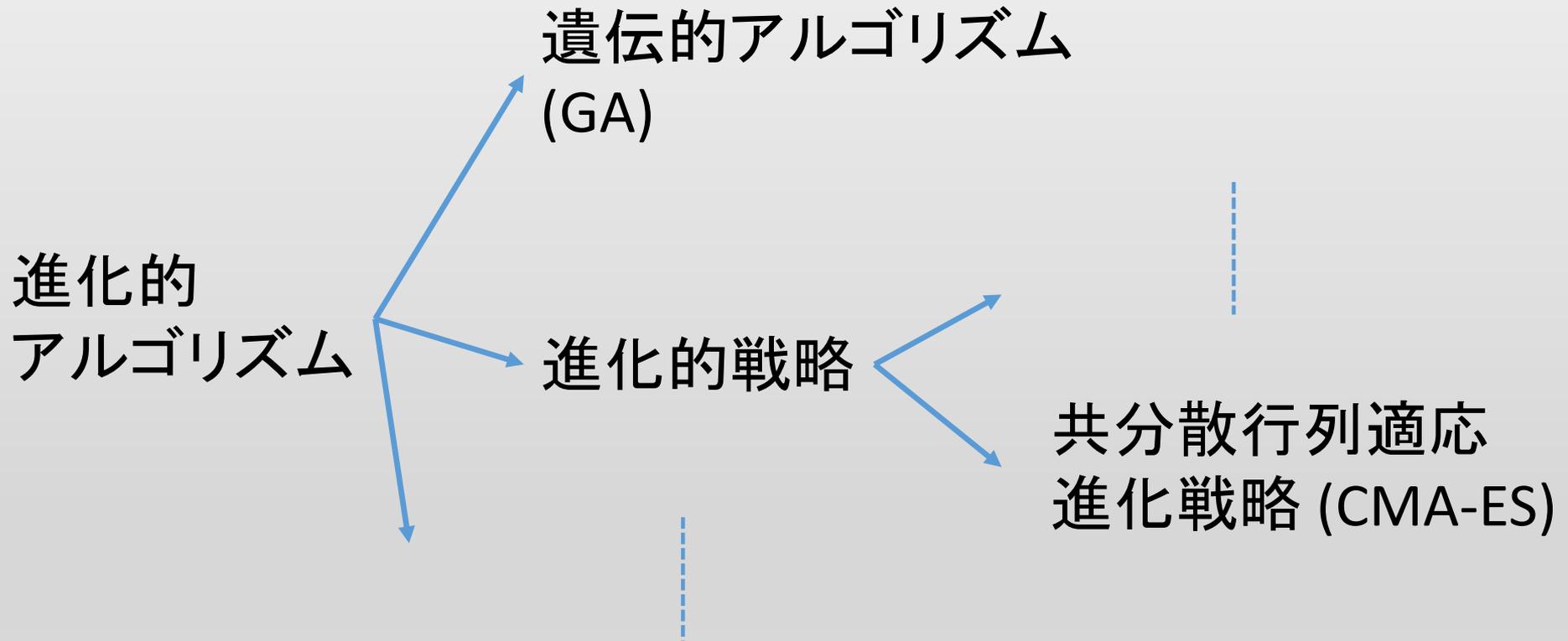
従来アプローチ
専門家によるチューニング



提案アプローチ
進化計算による自動チューニング



進化的アルゴリズムの種類



CMA-ESによる音声認識システムの自動最適化

① メタパラメタを遺伝子
(実数値ベクトル)に
エンコード

初期メタパラメタ

- 隠れ層の層数
- 各層のユニット数
- 学習係数
- etc.

② 遺伝子分布
をガウス分布
で表現

ガウス分布

③ サンプルングにより
遺伝子個体群を生成

x_1 x_2 x_3 x_4 x_N

繰り返す

評価結果

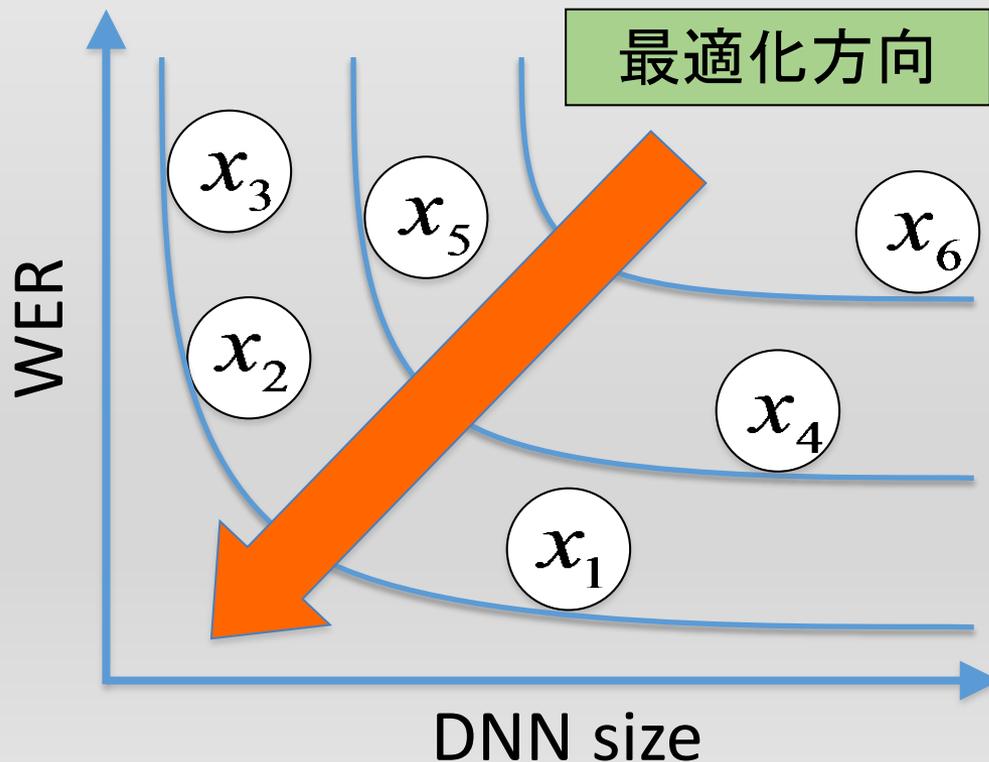
⑤ 評価結果を
もとに分布を更新

④ 各遺伝子から個体
を実体化・評価

パレート最適を用いた多目的最適化

認識性能と計算時間など、複数の目的関数を同時最適化

$$\left\{ \begin{array}{l} f_j(\mathbf{x}_k) \leq f_j(\mathbf{x}_{k'}) \quad \forall j = 1, \dots, J \\ \text{and} \\ f_j(\mathbf{x}_k) < f_j(\mathbf{x}_{k'}) \quad \exists j = 1, \dots, J \end{array} \right.$$



Ranking

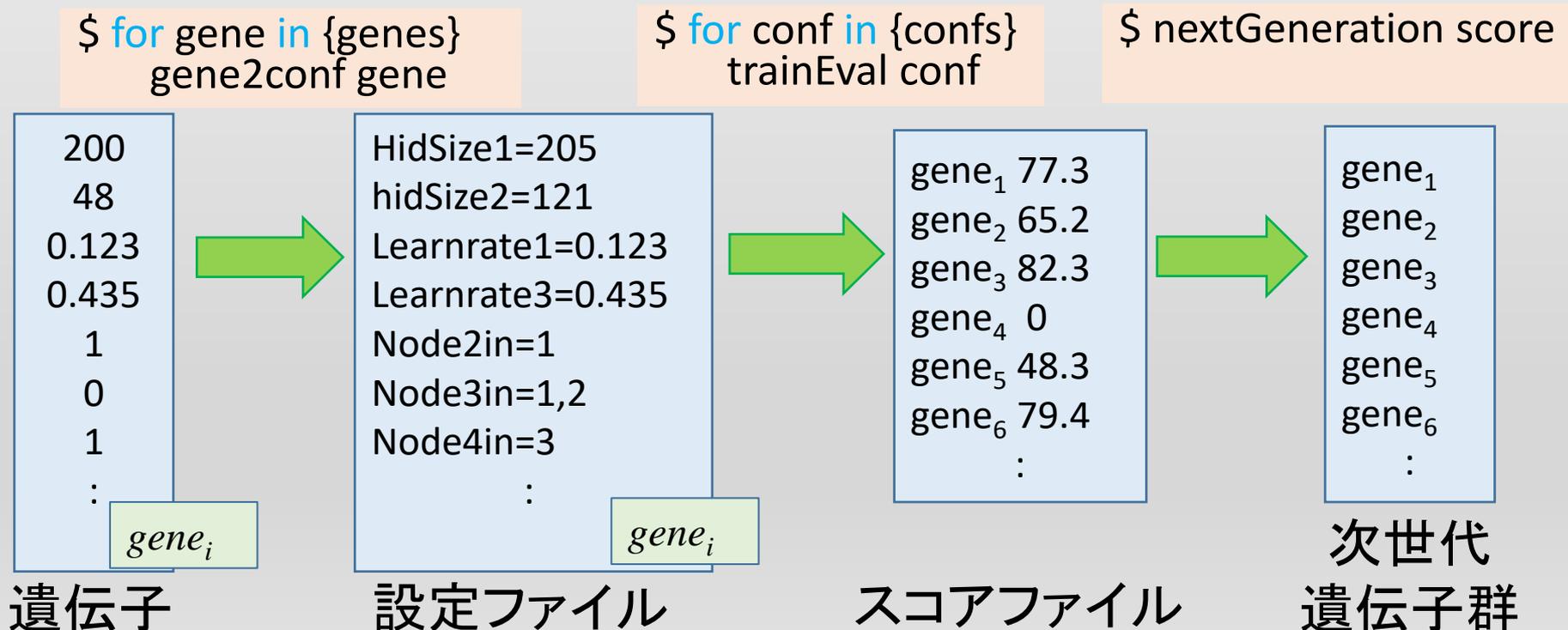
1. x_1
 x_2
 x_3
2. x_4
 x_5
3. x_6

目的関数として使用

進化フレームワークの実装

• トップレベルスクリプト

- gene2conf: 遺伝子を設定ファイルに変換(遺伝子の実体化)
- trainEval: モデル学習と評価(個体の評価)
- nextGeneration: ガウス分布の更新と次世代遺伝子群の生成



個体群の並列評価

1 個体の評価 = ニューラルネットワークの学習と評価

全個体の評価には長時間の計算が必要

しかし、世代内においては、各個体の学習・評価は
完全に独立して実行できる



1 個体に1GPUを割り当てて、
並列計算



DNN音響モデルの進化最適化

実験条件

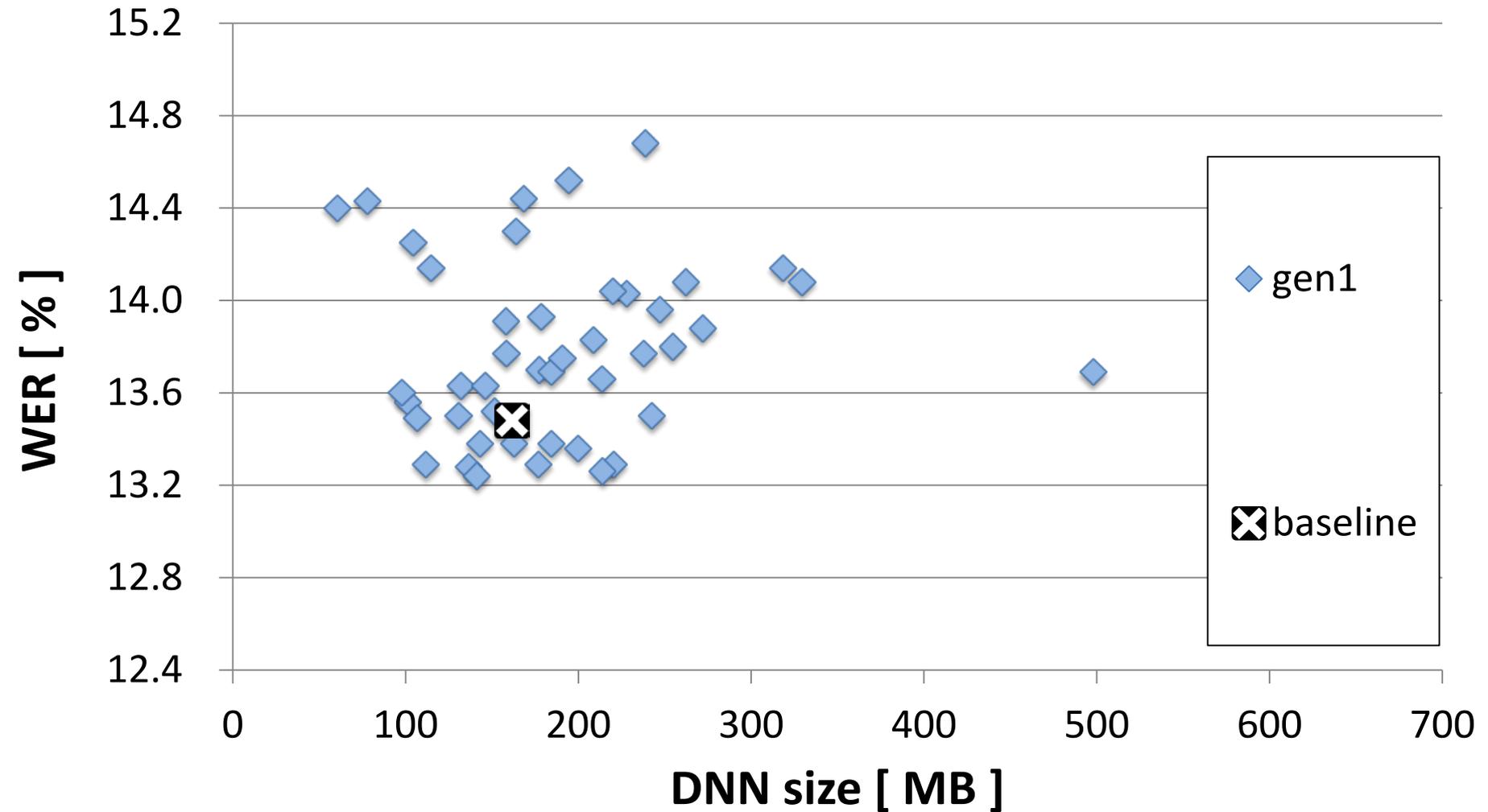
データベース	日本語話し言葉コーパス (CSJ)
認識システム構築	Kaldi Speech Recognition Toolkit
学習データ	240 時間
個体(並列)数 / 世代	44
使用計算機	TSUBAME 2.5 (GPU : Tesla K20X)
開発・評価セット	それぞれ2 時間(10講演)の音声データ

進化の評価尺度	単語正解精度 + DNNファイルサイズ
---------	---------------------

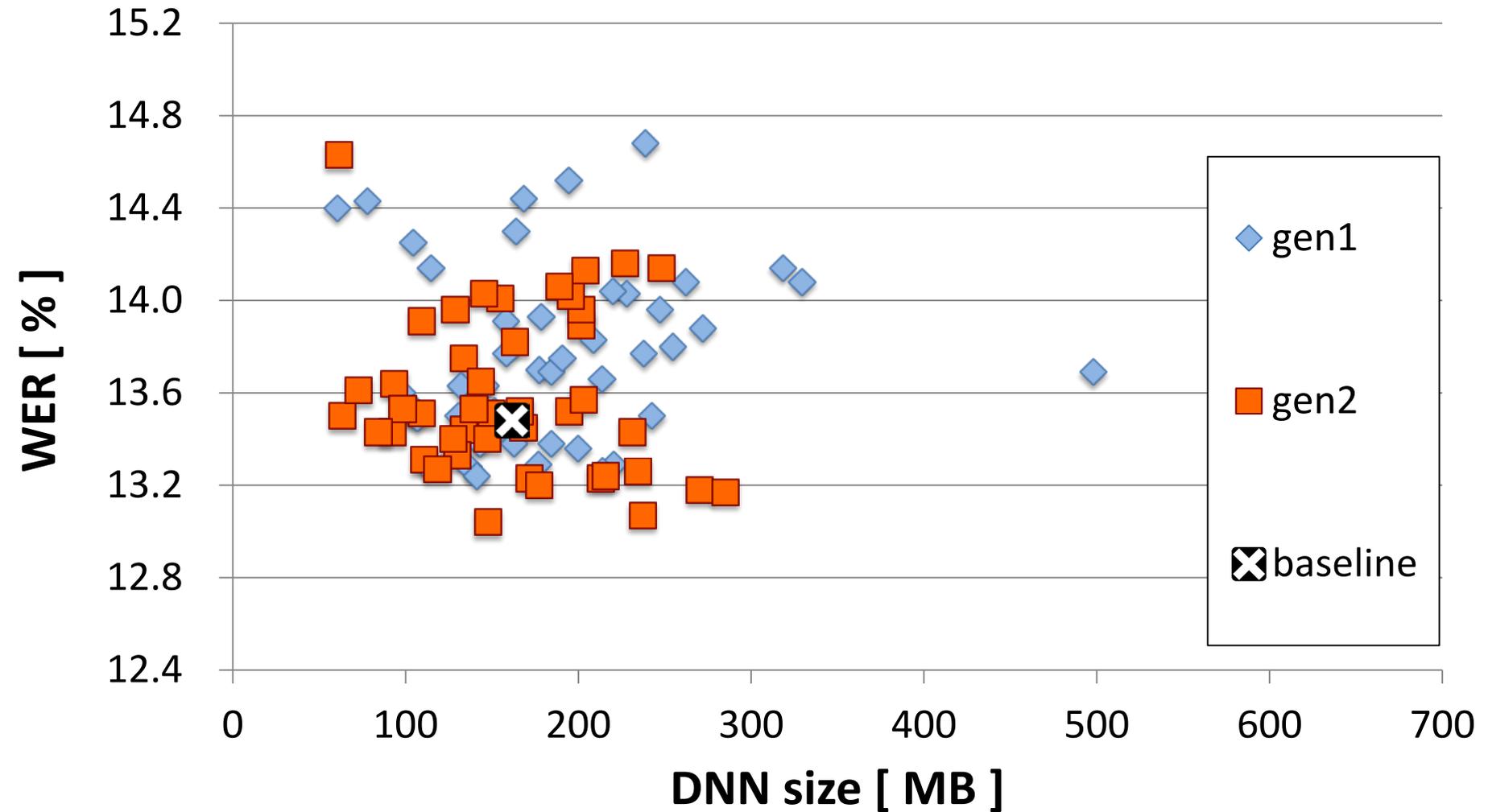
最適化対象のメタパラメタ

カテゴリ	メタパラメタ名	メタパラメタの詳細	初期値 x_{init}
特徴量	feat_type	特徴量の種類(MFCC,PLP,FBANK)	MFCC
DNN構造	splice	入力特徴量のフレーム幅	5
	nn_depth	中間層の数	6
	hid_dim	中間層1層あたりの素子数	2048
RBM学習 パラメタ	param_stddev_first	パラメタ初期化の分散値	0.1
	param_stddev		0.1
	rbm_lrate	学習率	0.4
	rbm_lrate_low		0.01
	rbm_l2penalty	L2正則化係数	0.0002
DNN学習 パラメタ	learn_rate	学習率	0.008
	momentum	モーメントム	0.00001

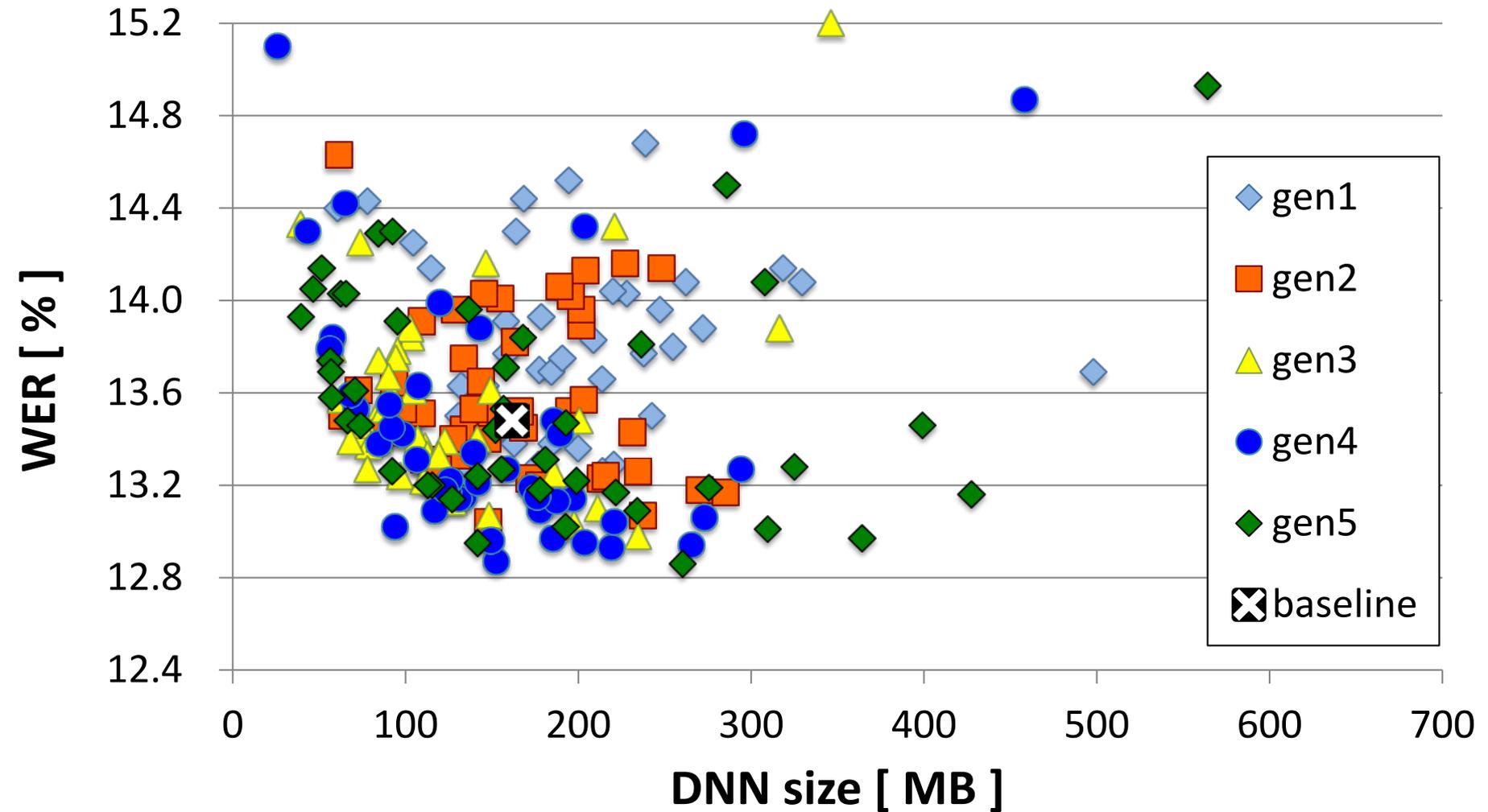
各世代の性能分布



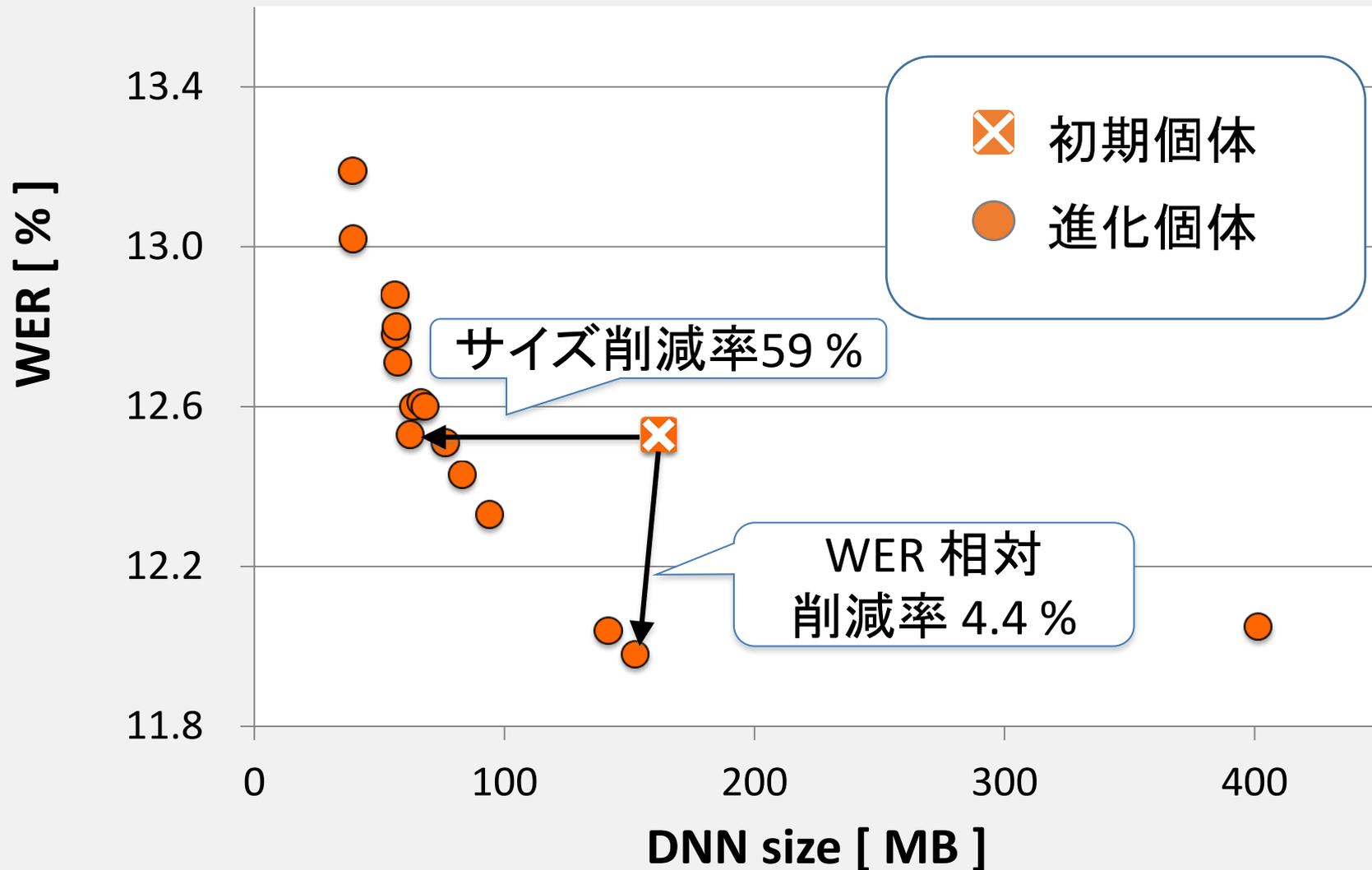
各世代の性能分布



各世代の性能分布



全個体によるパレートフロント



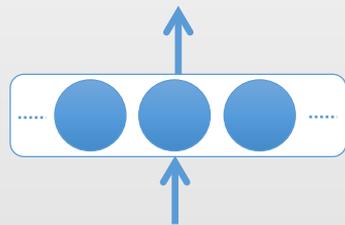
RNN言語モデルの進化最適化

最適化対象のメタパラメタ

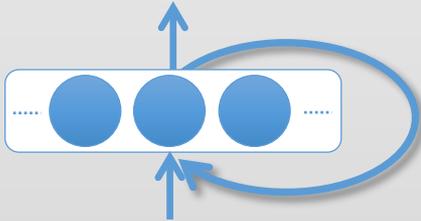
- ネットワーク構造
 - 隠れ層の数, 各層のユニット数
 - 語彙サイズ (= 入力層, 出力層のユニット数)
 - 各層のユニットタイプ (LSTM, RNN, FF)
- 各種調整値
 - 学習率の初期値, モーメントム
 - ミニバッチサイズ
 - 学習器の種類とそのメタパラメタetc.

隠れ層のユニットタイプ

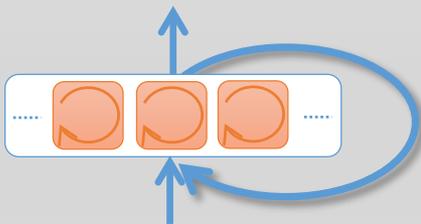
層毎のユニットタイプを可変にして進化により選択



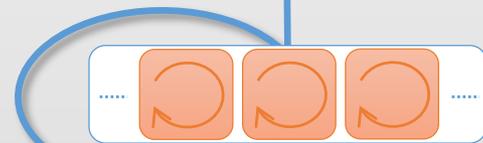
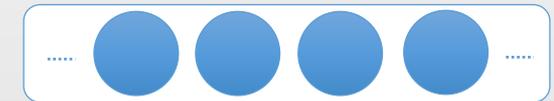
FF layer
最も単純な構造



RNN layer
再帰構造により過去の
文脈情報を利用可能



LSTM layer
RNN layer より長い
文脈情報が利用可能



LSTM
layer



FF
layer



RNN
layer



例

*LSTM: Long Short-Term Memory

実験条件

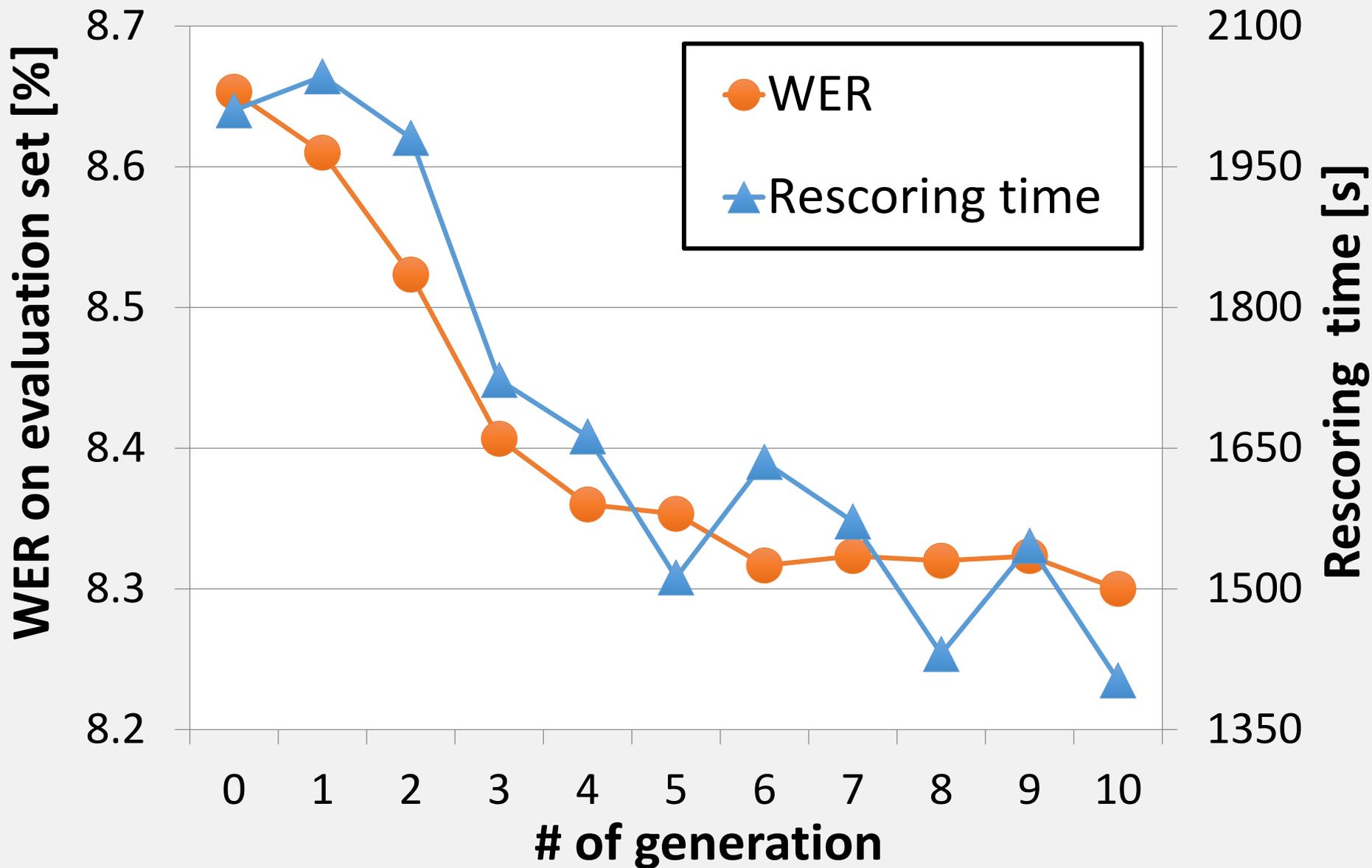
- 全般条件

データベース		日本語話し言葉コーパス (CSJ)
学習セット	音響モデル	520 時間
	言語モデル	750万 単語
評価セット		6 時間
認識デコーダ		Kaldi

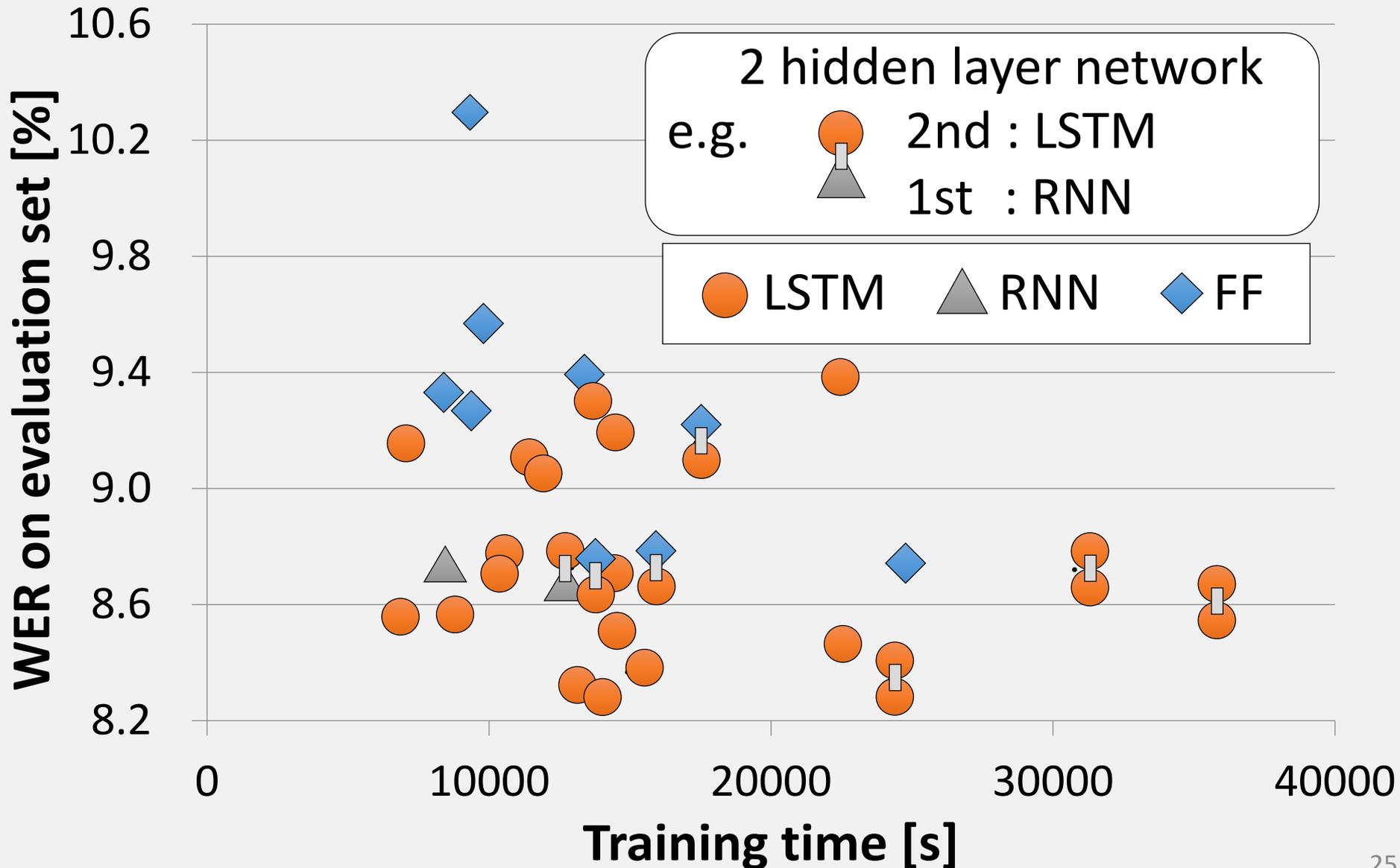
- 進化条件

開発セット	約 30 時間
世代数	10 世代
個体数	30 個体
評価尺度	単語正解精度 + RNN言語モデルの計算時間

進化結果



モデル構造と性能の分布 (10世代目)



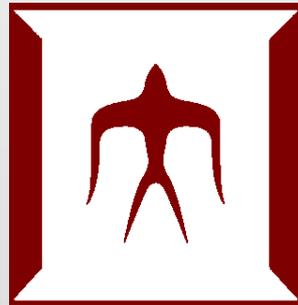
まとめと今後の課題

まとめ

- 進化的アルゴリズムにより、音声認識システムのシステム開発でボトルネックとなっているメタパラメタチューニングを自動化
- パレート最適を用いることで、単語誤り率WERと計算時間など複数の尺度を同時最適化
- 開発した高精度日本語音声認識システムは、Kaldiパッケージに収録され一般利用が可能 <http://kaldi-asr.org/>

今後の課題

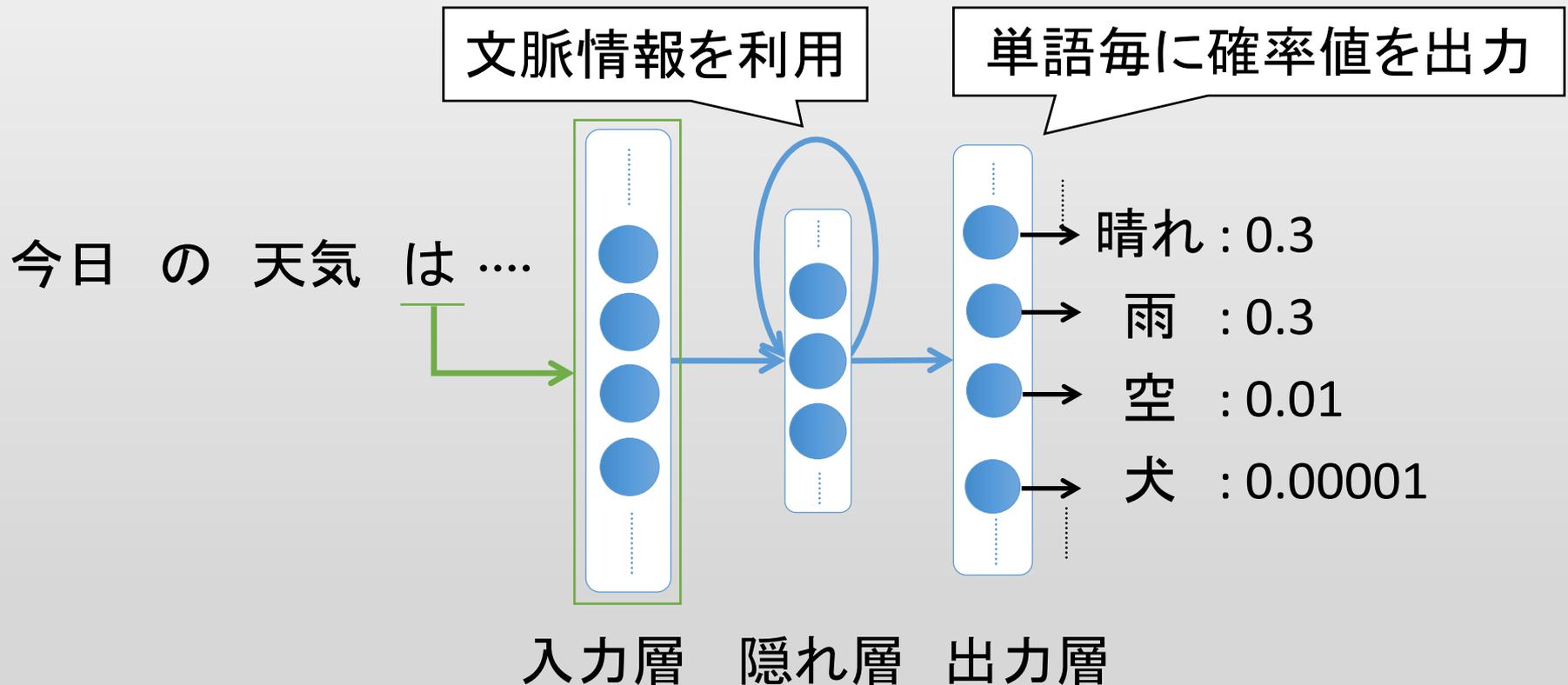
- より複雑なネットワーク構造や多様なユニットタイプを対象とした進化実験
- 進化効率の向上と超大規模データへのスケールアップ



END

補足資料

RNN言語モデルの基本構造



LSTM layer

- RNN layer の各ユニットをLSTMブロックに置換
- RNN layer より長い文脈情報を利用可能

