

日本音響学会2010年秋季研究発表会

音声 F_0 パターン生成過程の確率モデル

○ 亀岡弘和
ルルージョナトン
大石康智

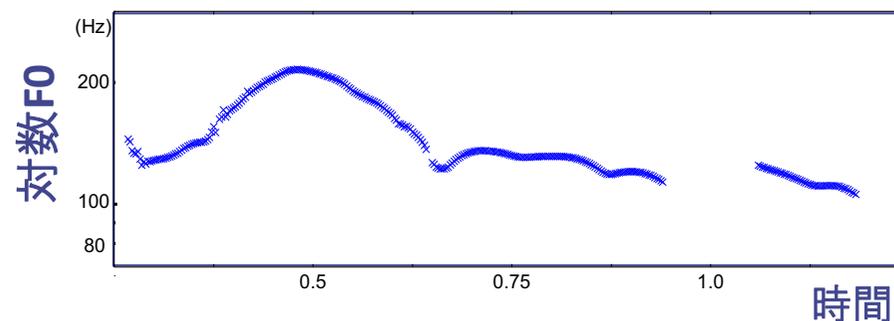
日本電信電話株式会社
NTT コミュニケーション科学基礎研究所

音声の基本周波数(F0)パターンについて

◆F0パターンとは？

→ 声帯振動の周期の時間変化のパターン

- ◆ 甲状軟骨の運動
⇒ 声帯の長さが微小変化 ⇒ 固有振動数が変化



◆音声対話において果たす役割は？

→ 音声の韻律的特徴を表現

- ◆ 構文、意味、感情、焦点

◆音声情報処理における有用性

→ F0パターンの自然性は音声らしさを表す有用な指標の一つ

「F0パターンが生成過程の物理的制約をどれだけ満たしているか」

→ 音声らしさを測るための客観的な規準に！

F0パターンのモデル化・分析は
様々な音声アプリケーションに有益なはず → 本研究の動機

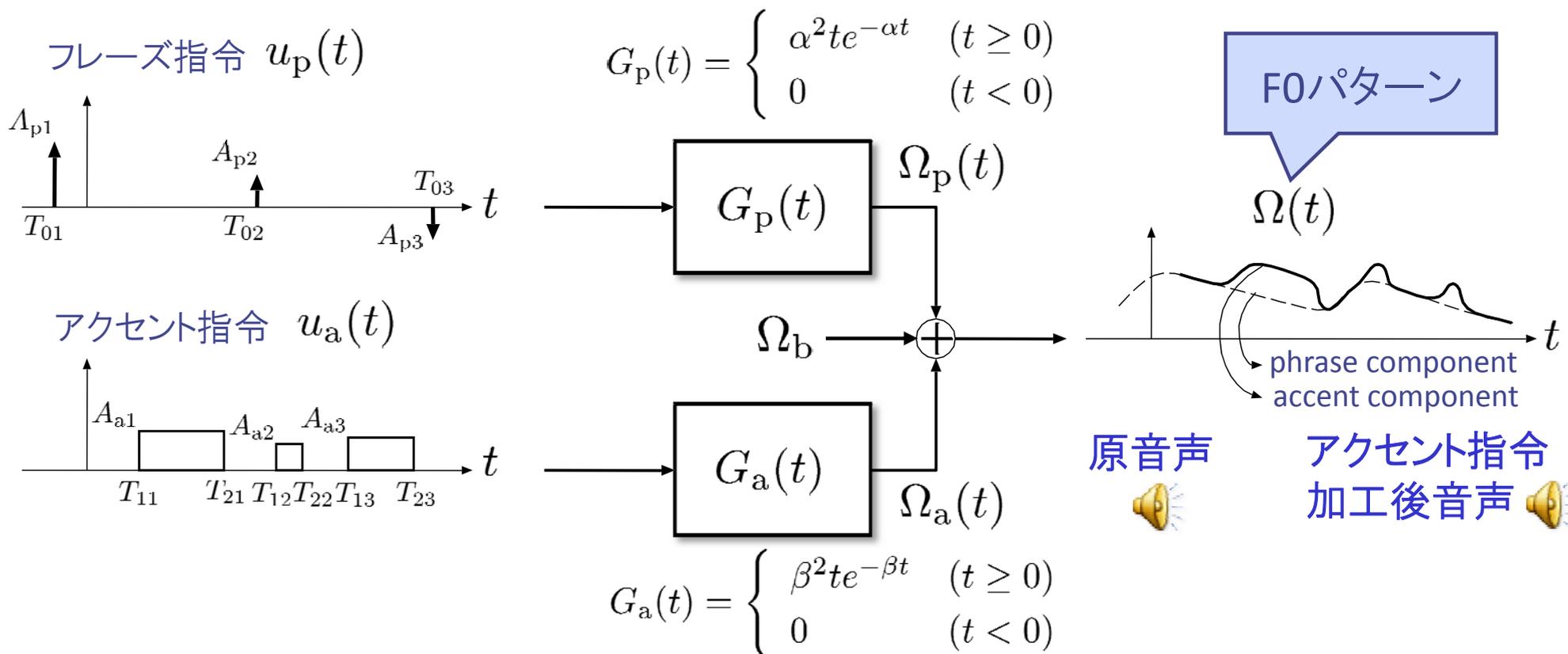
藤崎のF0パターン生成過程モデル(藤崎モデル) [藤崎1988]

[仮定1] 甲状軟骨の運動は二種類(平行移動と回転)

[仮定2] 各運動を臨界制動の二次線形系の応答として表現

[仮定3] 各運動に伴う声帯の伸びの変化の合計が対数F0の変化に比例

フレーズ成分 $\Omega_p(t)$: 平行移動による変位に比例するF0パターン成分
 アクセント成分 $\Omega_a(t)$: 回転運動による変位に比例するF0パターン成分
 ベースライン成分 Ω_b : 声帯の物理的性質によって決まる定数



本研究の目的

(1) 藤崎モデルの離散時間表現を統計モデル化 (2)

- 統計的手法を駆使した強力なパラメータ推定法の枠組を確立
- 統計モデルに基づく音声処理問題(合成、分析、分離、強調)にスムーズに組み込める基盤を構築

発表アウトライン

◆モデル化の手順

(1) 後退差分近似による連続時間藤崎モデルの離散化

(2) フレーズ指令 & アクセント指令の確率モデル化

→ (1)と(2)からF0パターンの確率密度関数を導出

◆提案モデルの使用方法

◆パラメータ推定アルゴリズムの導出

本研究の目的

(1) 藤崎モデルの離散時間表現を統計モデル化 (2)

- 統計的手法を駆使した強力なパラメータ推定法の枠組を確立
- 統計モデルに基づく音声処理問題(合成、分析、分離、強調)にスムーズに組み込める基盤を構築

発表アウトライン

◆モデル化の手順

(1) 後退差分近似による連続時間藤崎モデルの離散化

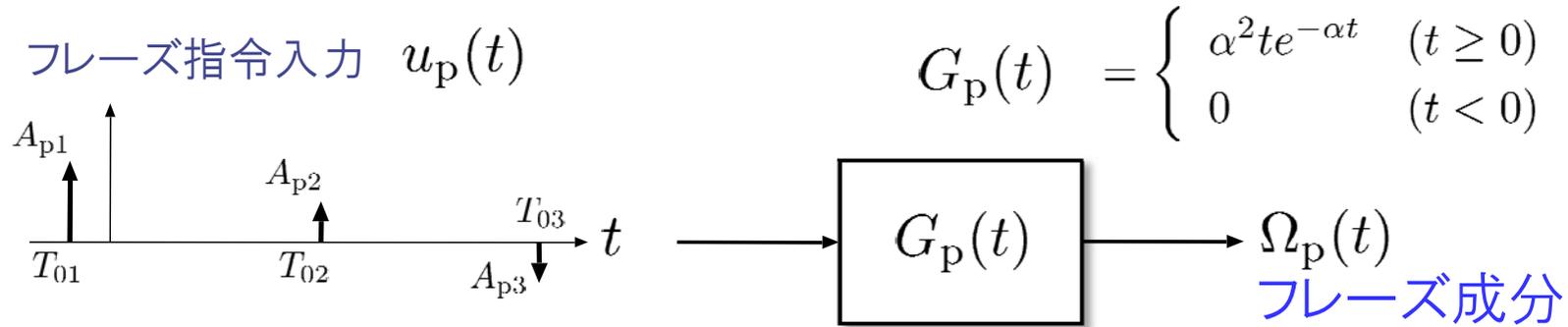
(2) フレーズ指令 & アクセント指令の確率モデル化

→ (1)と(2)からF0パターンの確率密度関数を導出

◆提案モデルの使用方法

◆パラメータ推定アルゴリズムの導出

連続時間系から離散時間系へ(フーズ制御系)



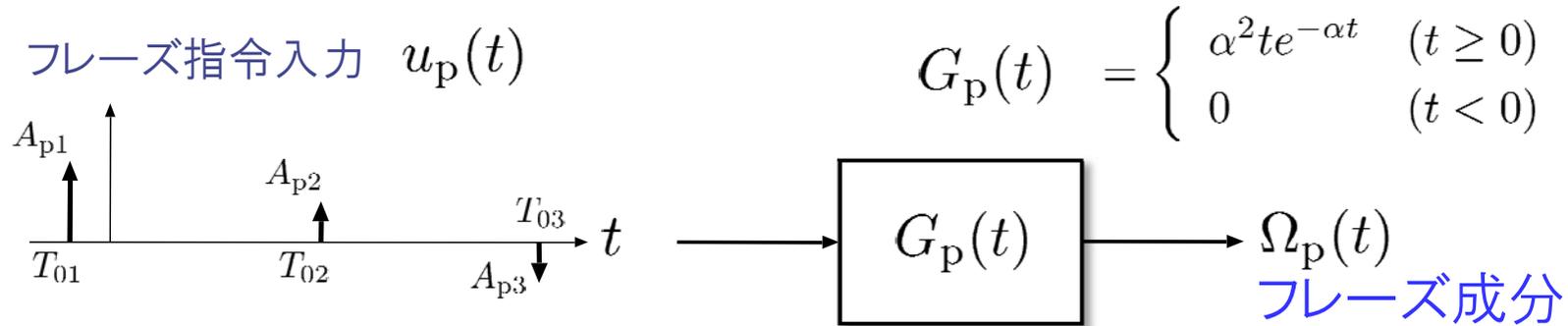
- フーズ制御系のラプラス変換 $\alpha^2 / (s + \alpha)^2$
- 逆システムのラプラス変換 $\rightarrow (s + \alpha)^2 / \alpha^2$
- 後退差分近似 $\rightarrow s = \frac{1 - z^{-1}}{t_0}$ → 離散時間表現のサンプリング周期
- 逆システムの伝達関数 $\frac{z^{-2} - 2(\alpha t_0 + 1)z^{-1} + (\alpha t_0 + 1)^2}{\alpha^2 t_0^2}$

➡ $u_p[l] = g_0^P \Omega_p[l] + g_1^P \Omega_p[l - 1] + g_2^P \Omega_p[l - 2]$

$$g_0^P = (\psi - 1)^2, \quad g_1^P = -2\psi(\psi - 1), \quad g_2^P = \psi^2, \quad \psi = 1 + 1/(\alpha t_0)$$

パラメータ ψ で特性が決まる拘束つき全極モデル

連続時間系から離散時間系へ(フレーズ制御系)



■ フレーズ制御系のラプラス変換 $\alpha^2 / (s + \alpha)^2$

■ 逆システムのラプラス変換 $\rightarrow (s + \alpha)^2 / \alpha^2$

■ 後退差分近似

■ 逆システムの

アクセント制御系も同形
なので以上と同様！

表現のサンプリング周期

$$\frac{1}{1 + (\alpha t_0 + 1)^2}$$

➡ $u_p[l] = g_0^P \Omega_p[l] + g_1^P \Omega_p[l - 1] + g_2^P \Omega_p[l - 2]$

$$g_0^P = (\psi - 1)^2, \quad g_1^P = -2\psi(\psi - 1), \quad g_2^P = \psi^2, \quad \psi = 1 + 1/(\alpha t_0)$$

パラメータ ψ で特性が決まる 拘束つき全極モデル

本研究の目的

(1) 藤崎モデルの離散時間表現を統計モデル化 (2)

- 統計的手法を駆使した強力なパラメータ推定法の枠組を確立
- 統計モデルに基づく音声処理問題(合成、分析、分離、強調)にスムーズに組み込める基盤を構築

発表アウトライン

◆モデル化の手順

(1) 後退差分近似による連続時間藤崎モデルの離散化

(2) フレーズ指令 & アクセント指令の確率モデル化

→ (1)と(2)からF0パターンの確率密度関数を導出

◆提案モデルの使用方法

◆パラメータ推定アルゴリズムの導出

本研究の目的

(1) 藤崎モデルの離散時間表現を統計モデル化 (2)

- 統計的手法を駆使した強力なパラメータ推定法の枠組を確立
- 統計モデルに基づく音声処理問題(合成、分析、分離、強調)にスムーズに組み込める基盤を構築

発表アウトライン

◆モデル化の手順

(1) 後退差分近似による連続時間藤崎モデルの離散化

(2) フレーズ指令 & アクセント指令の確率モデル化

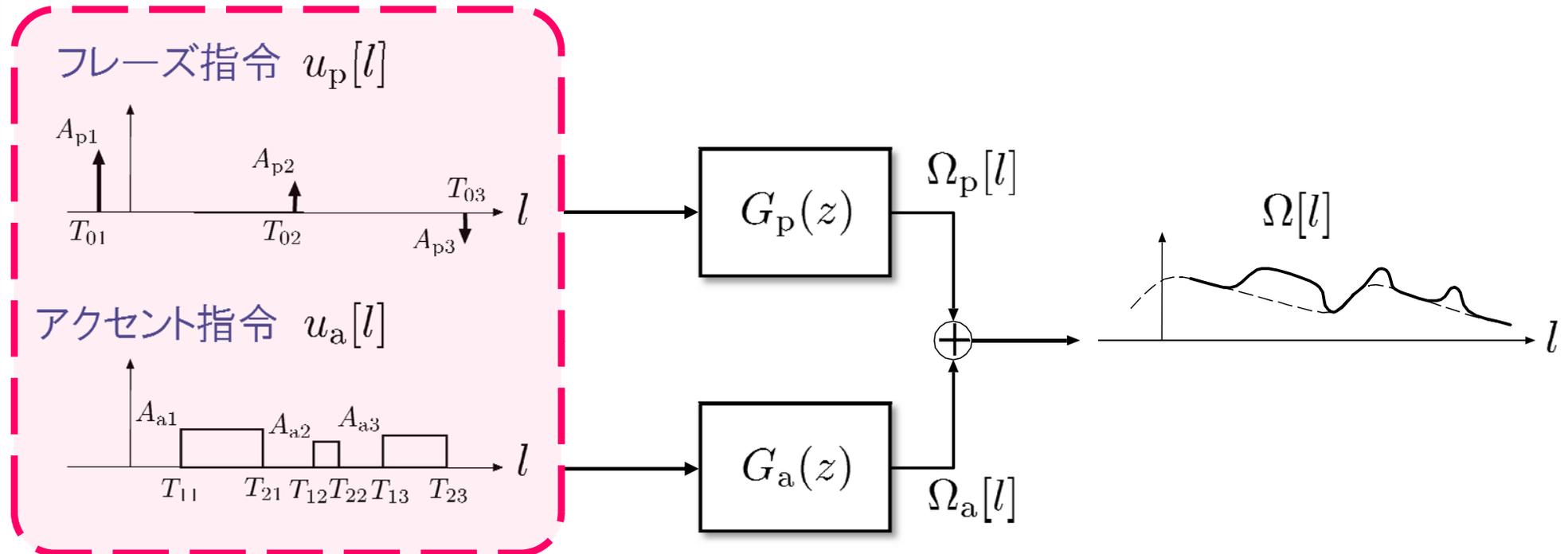
→ (1)と(2)からF0パターンの確率密度関数を導出

◆提案モデルの使用方法

◆パラメータ推定アルゴリズムの導出

フレーズ指令 & アクセント指令の確率モデル化

◆ 両指令は無制約で良い訳ではない

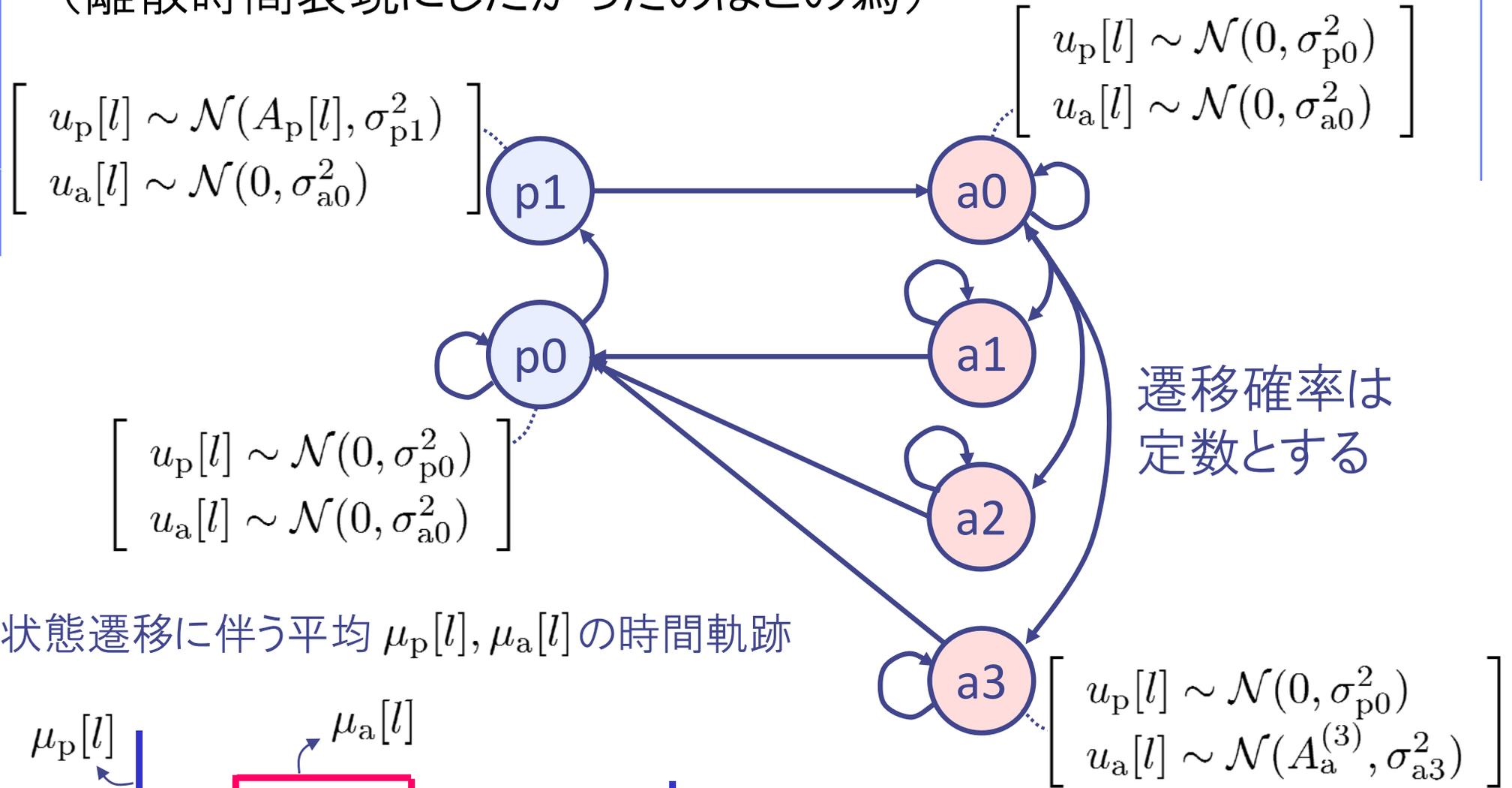


◆ 藤崎モデルにおける基本ルール

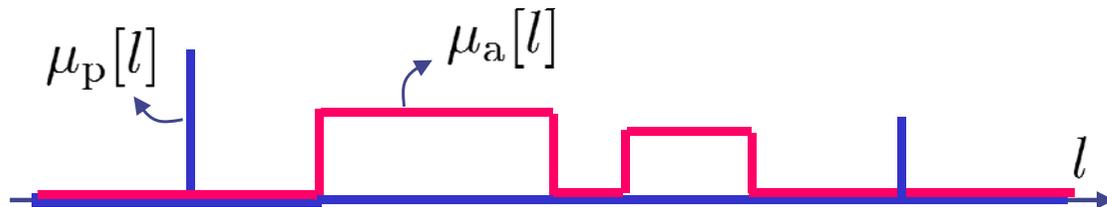
- フレーズ指令はパルス波
- アクセント指令は矩形波
- 指令は同時刻に2つ以上生起しない

フレーズ指令 & アクセント指令の確率モデル化

◆ 経路制限付き隠れマルコフモデルでモデル化
 (離散時間表現にしたかったのはこの為)

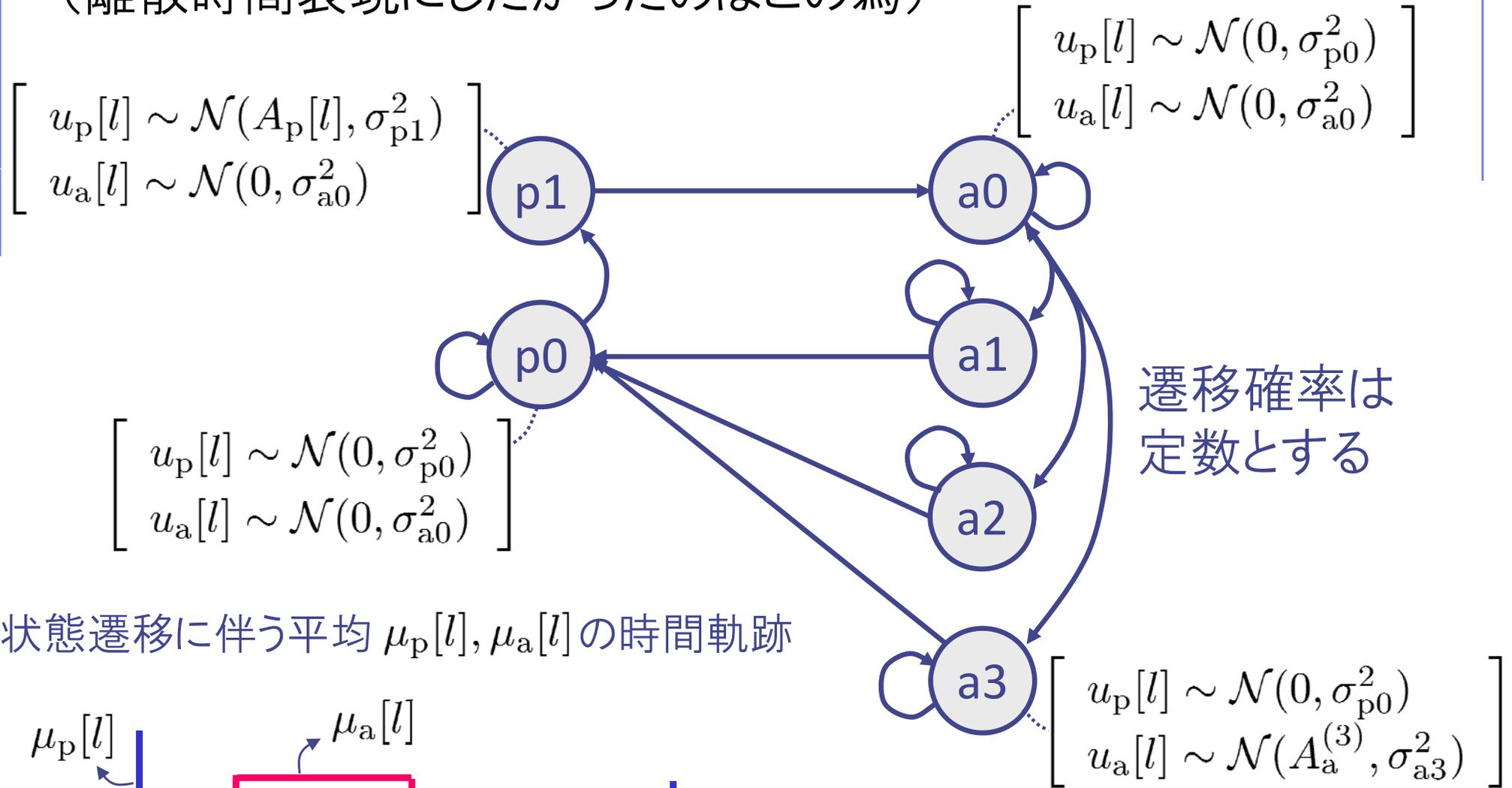


状態遷移に伴う平均 $\mu_p[l], \mu_a[l]$ の時間軌跡



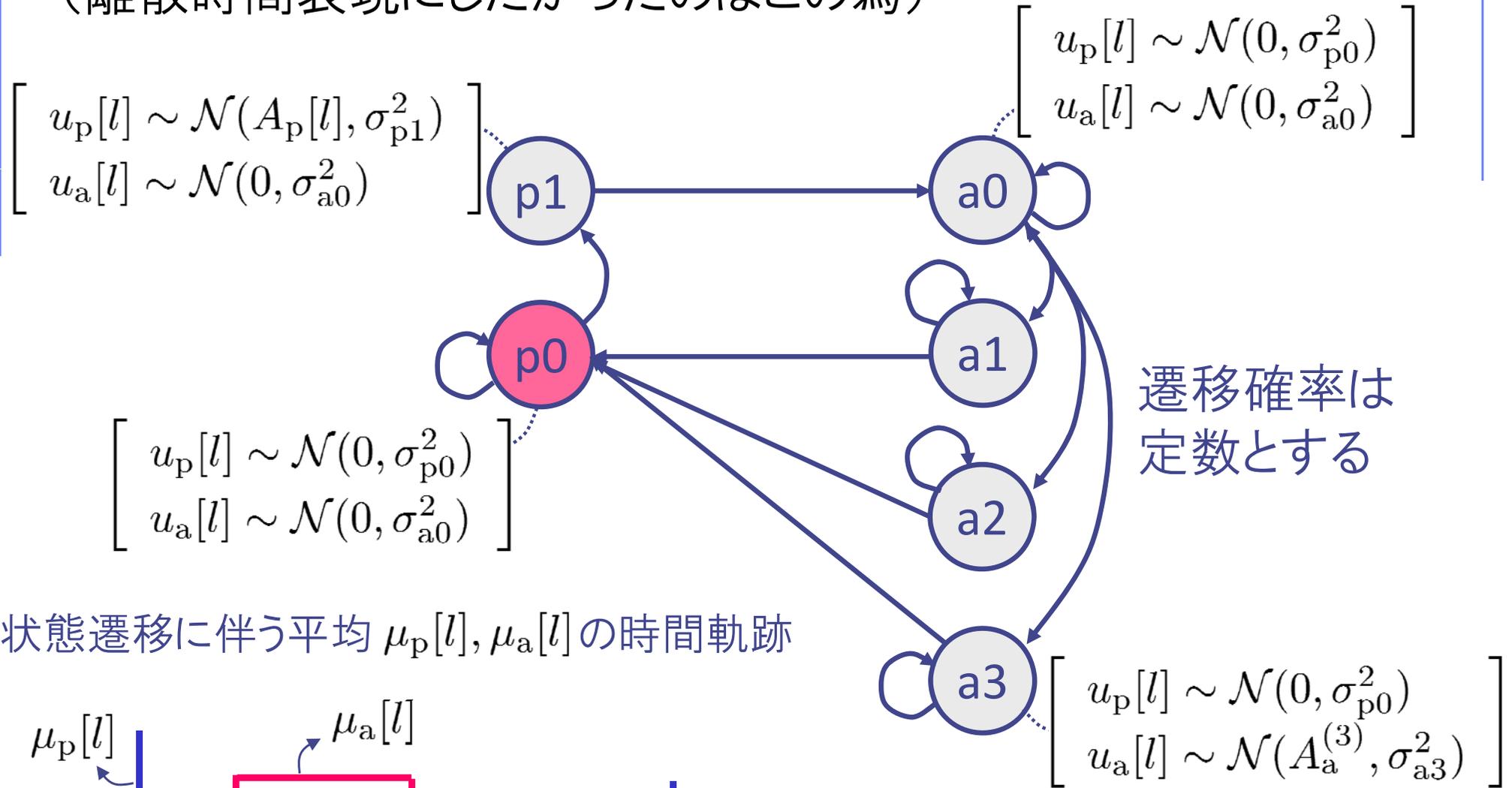
フレーズ指令 & アクセント指令の確率モデル化

◆ 経路制限付き隠れマルコフモデルでモデル化
 (離散時間表現にしたかったのはこの為)

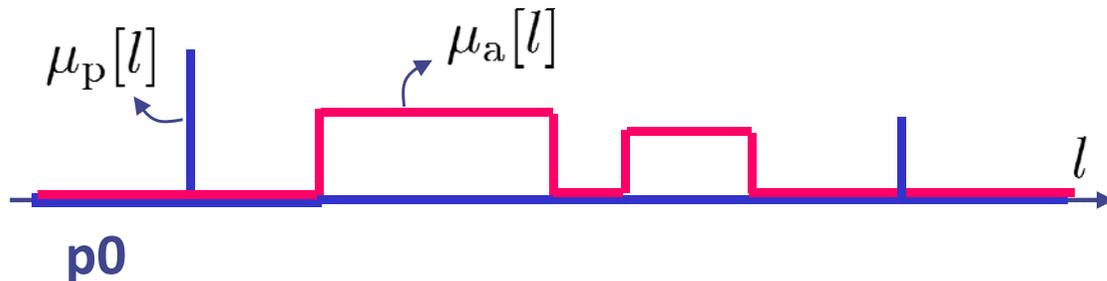


フレーズ指令 & アクセント指令の確率モデル化

◆ 経路制限付き隠れマルコフモデルでモデル化
 (離散時間表現にしたかったのはこの為)

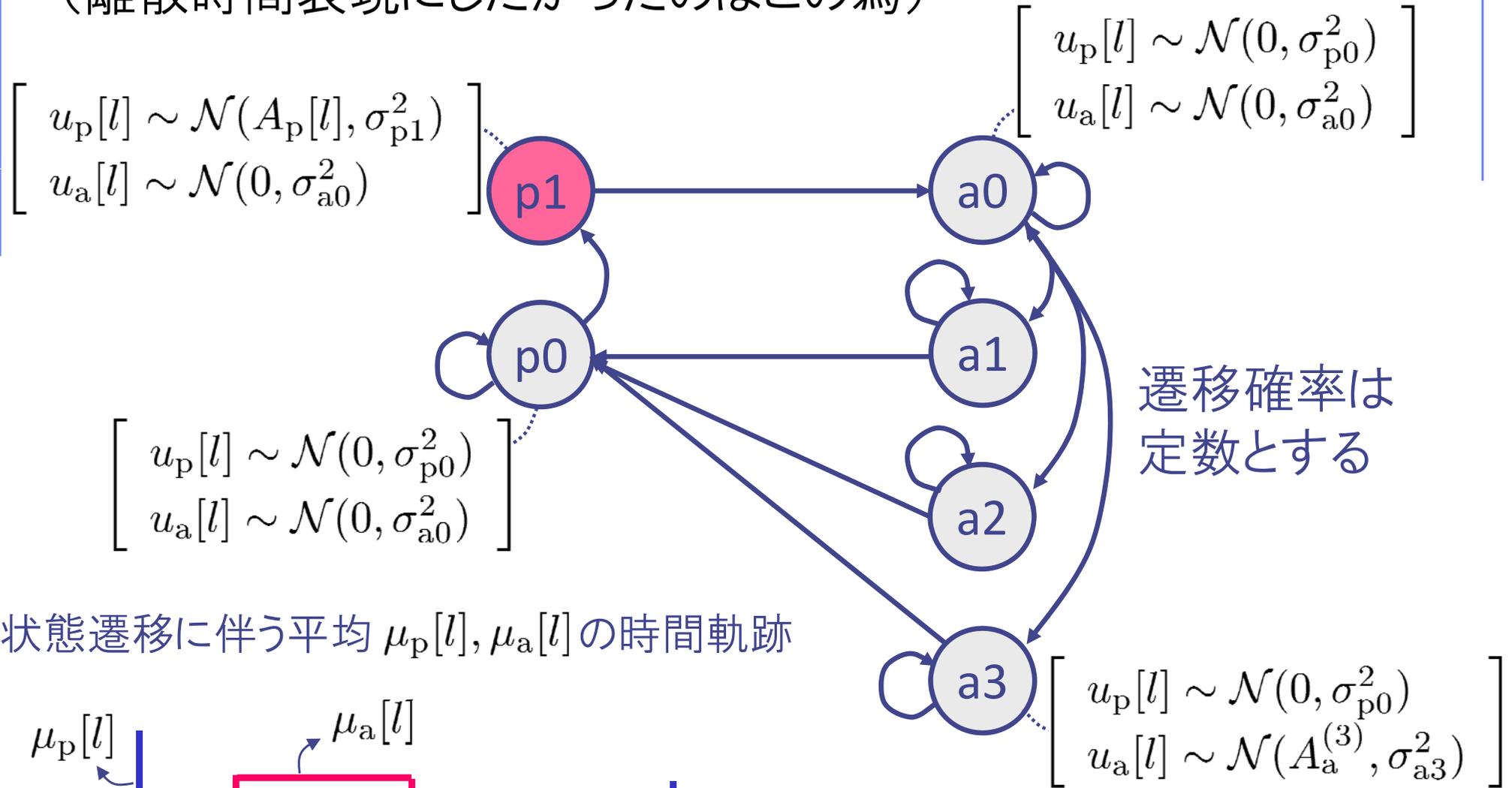


状態遷移に伴う平均 $\mu_p[l], \mu_a[l]$ の時間軌跡

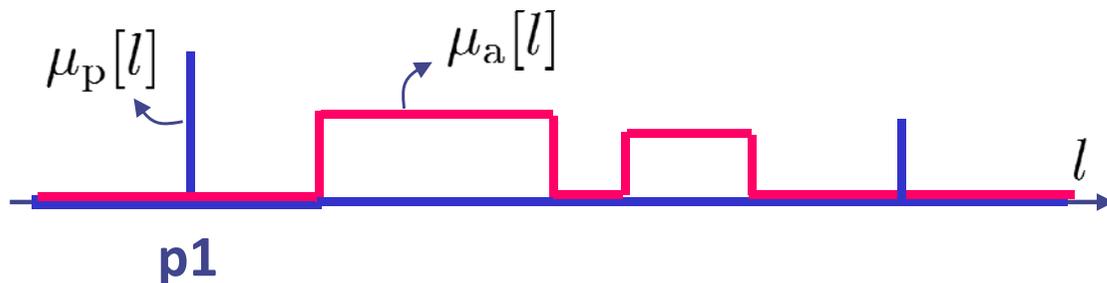


フレーズ指令 & アクセント指令の確率モデル化

◆ 経路制限付き隠れマルコフモデルでモデル化
 (離散時間表現にしたかったのはこの為)

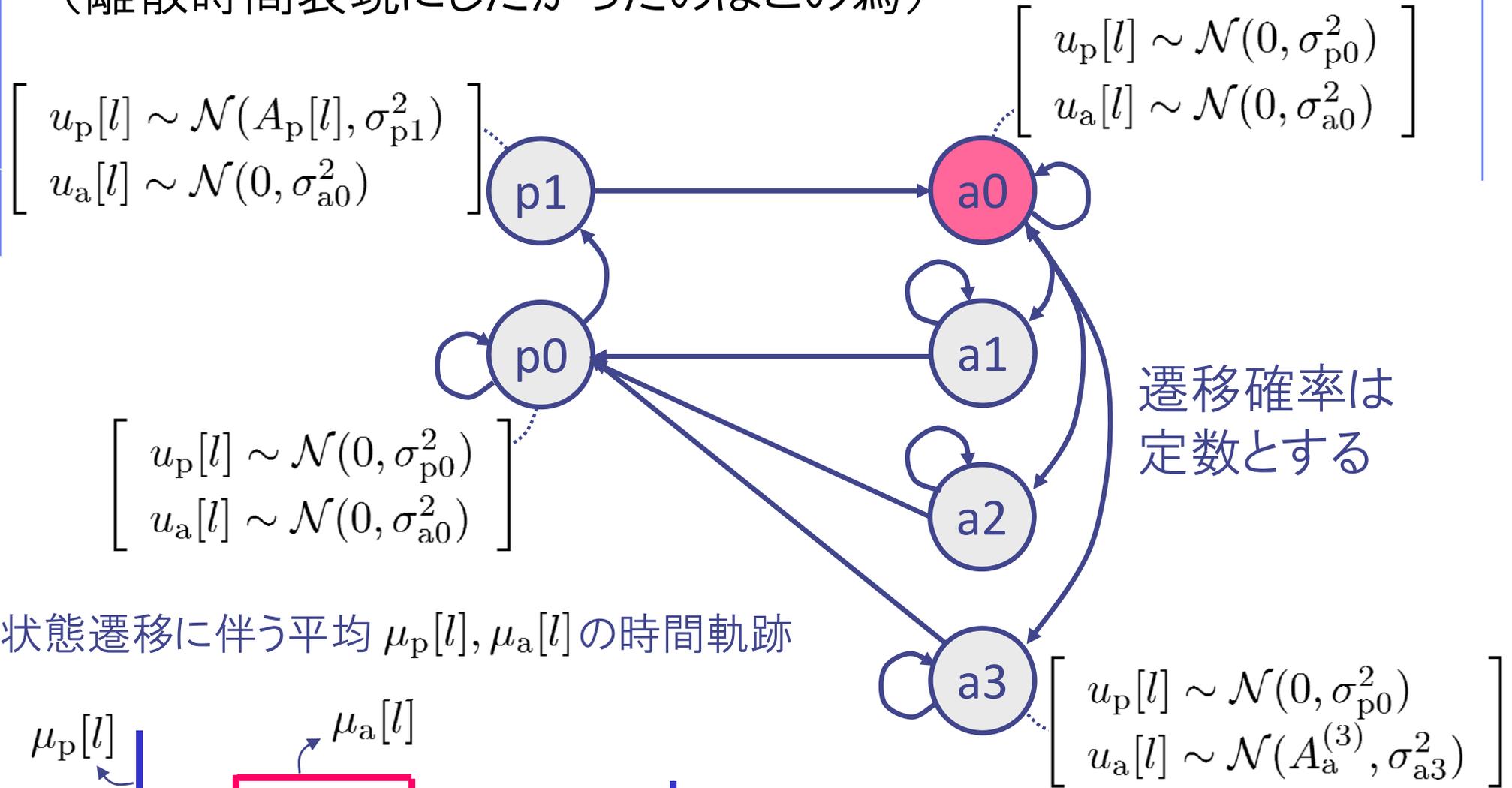


状態遷移に伴う平均 $\mu_p[l], \mu_a[l]$ の時間軌跡

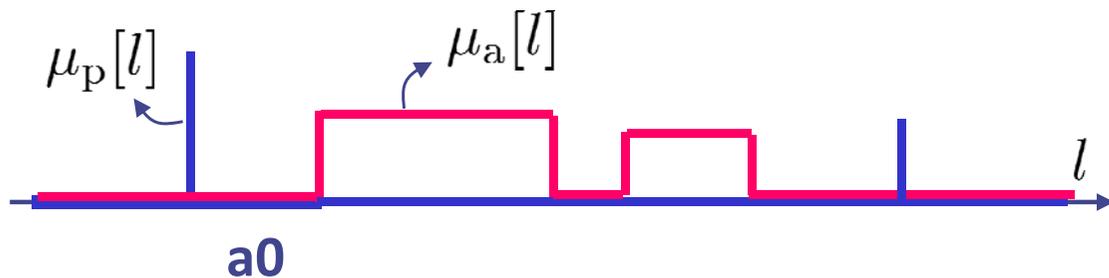


フレーズ指令 & アクセント指令の確率モデル化

◆ 経路制限付き隠れマルコフモデルでモデル化
 (離散時間表現にしたかったのはこの為)

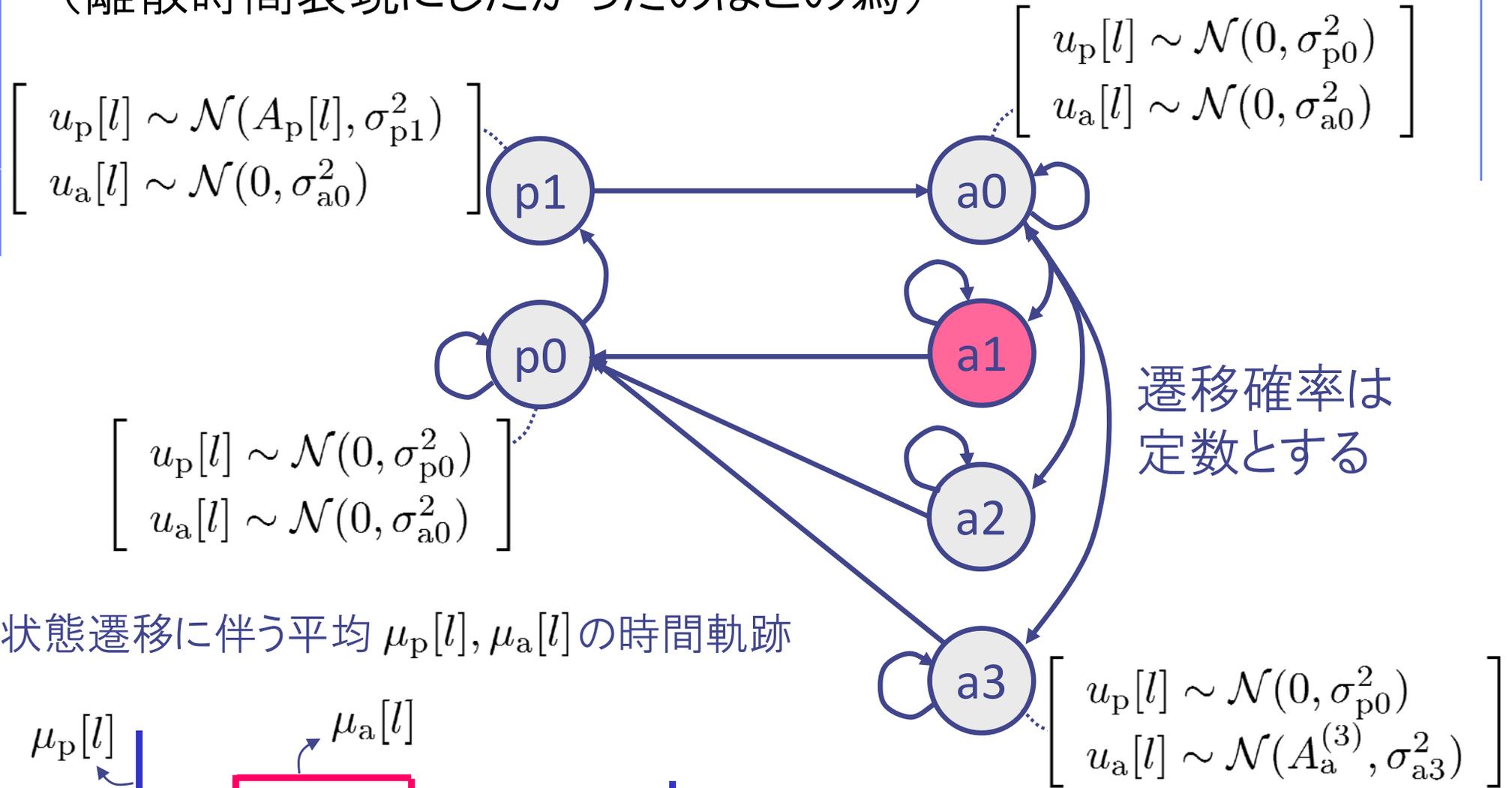


状態遷移に伴う平均 $\mu_p[l], \mu_a[l]$ の時間軌跡

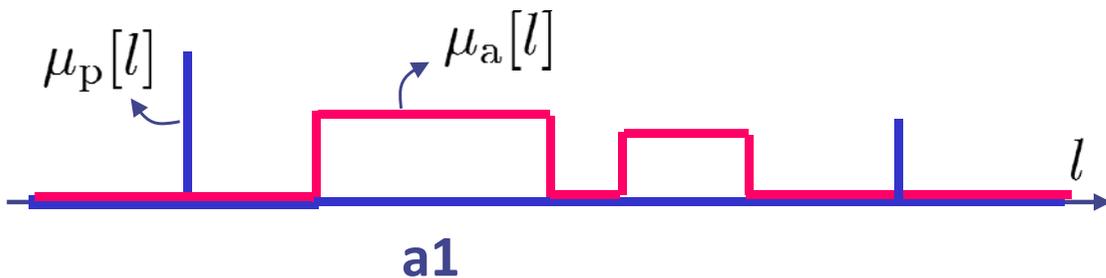


フレーズ指令 & アクセント指令の確率モデル化

◆ 経路制限付き隠れマルコフモデルでモデル化
 (離散時間表現にしたかったのはこの為)

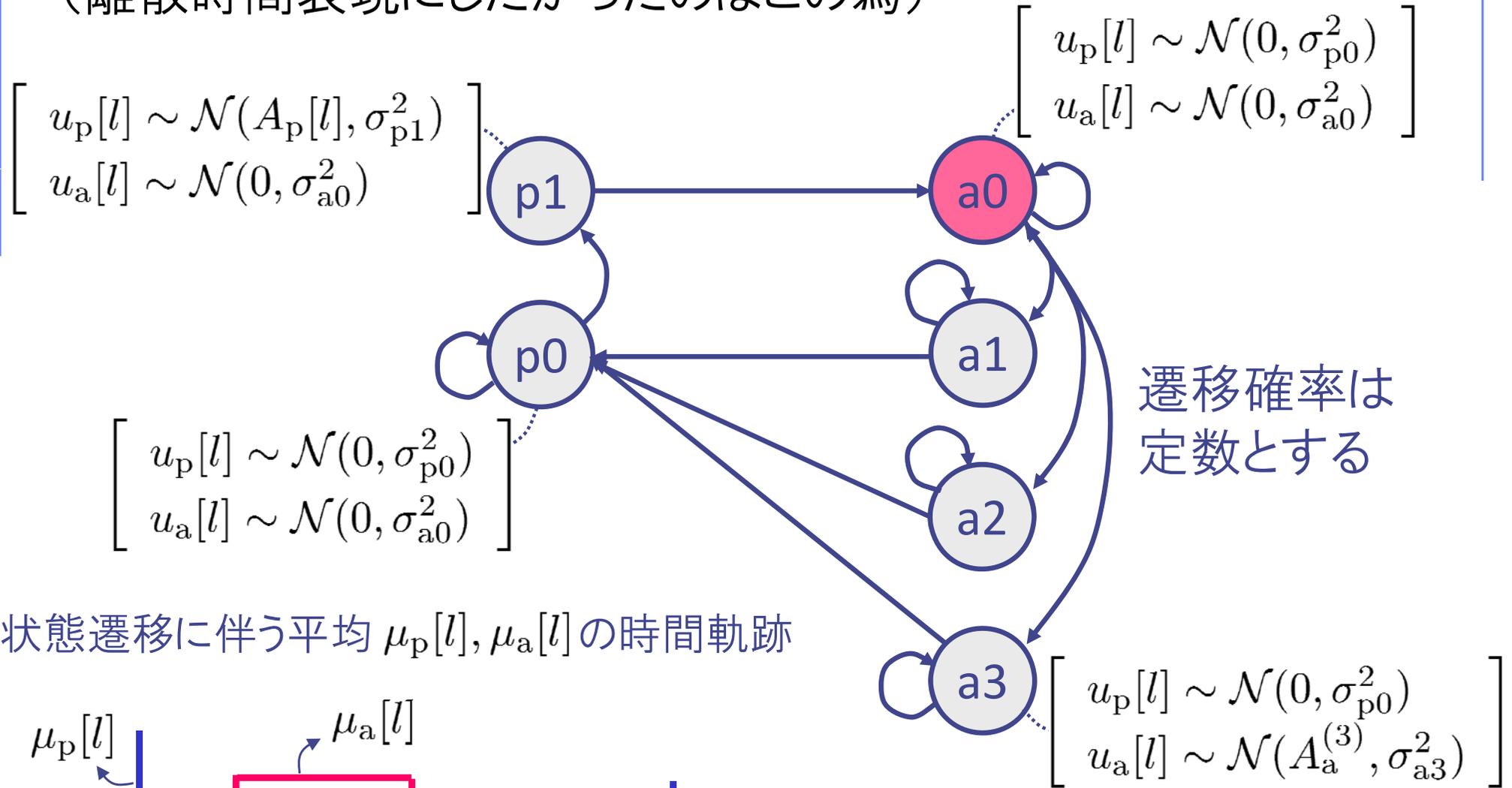


状態遷移に伴う平均 $\mu_p[l], \mu_a[l]$ の時間軌跡

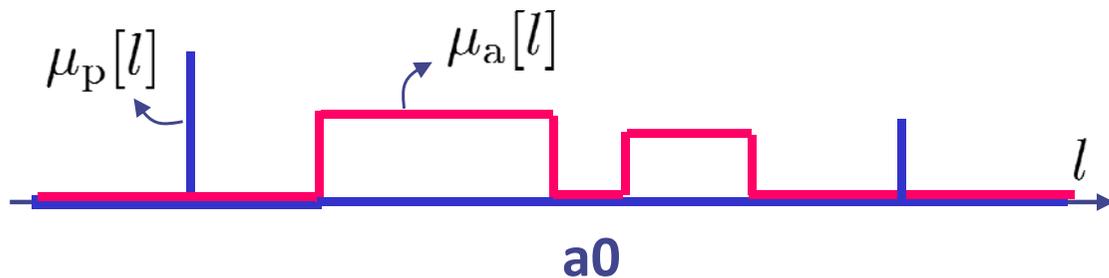


フレーズ指令 & アクセント指令の確率モデル化

◆ 経路制限付き隠れマルコフモデルでモデル化
 (離散時間表現にしたかったのはこの為)

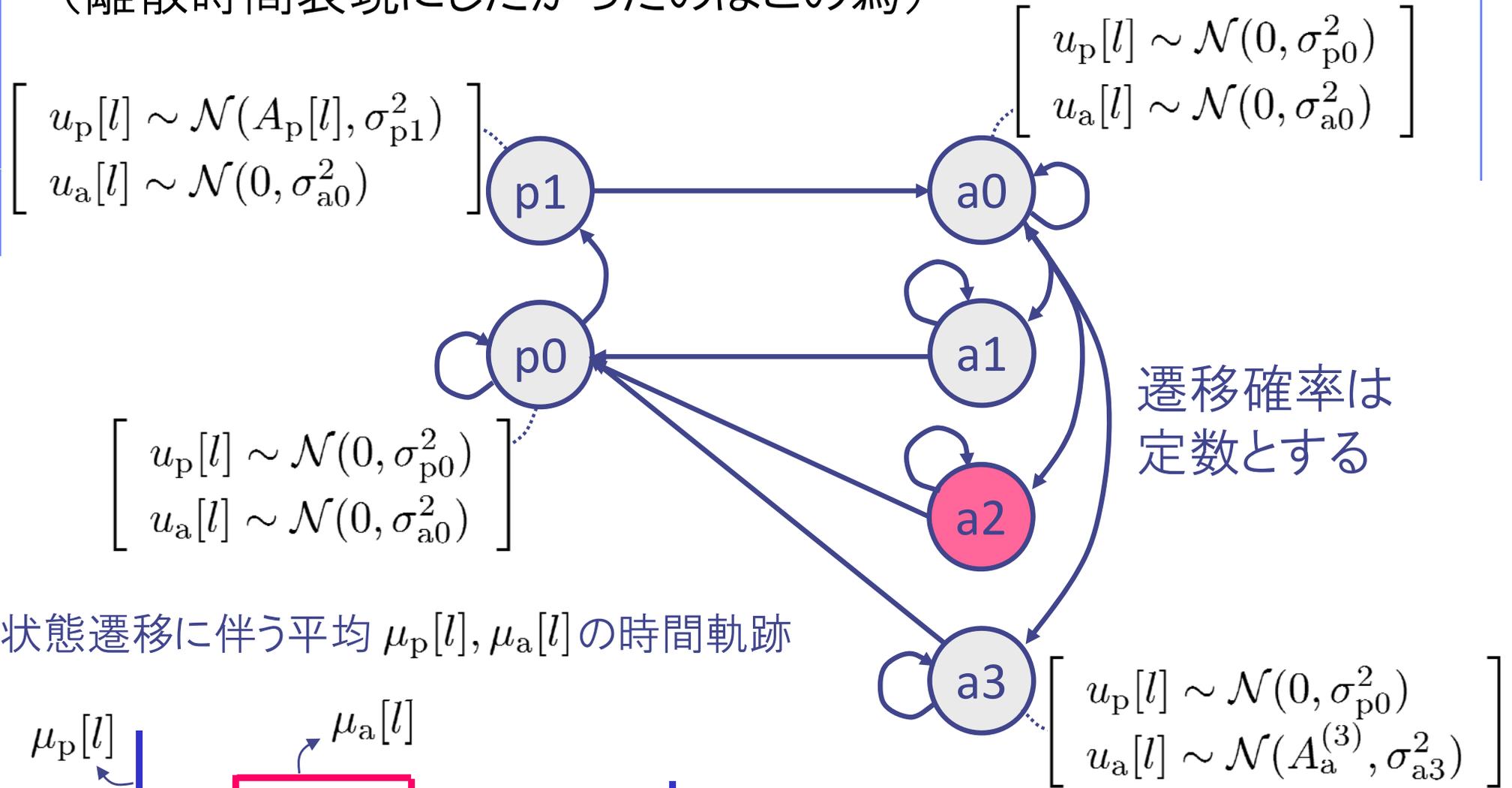


状態遷移に伴う平均 $\mu_p[l], \mu_a[l]$ の時間軌跡

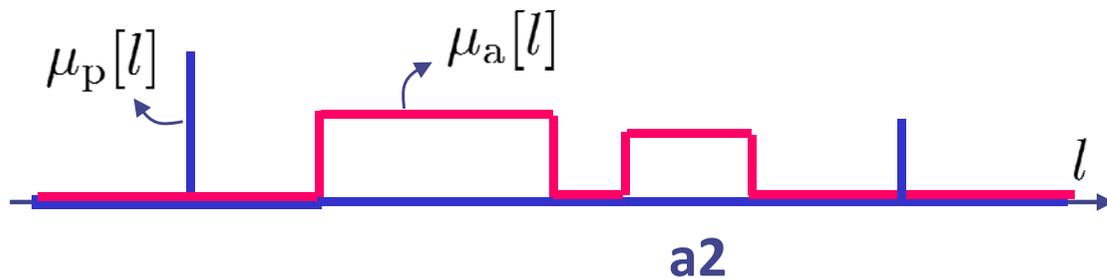


フレーズ指令 & アクセント指令の確率モデル化

◆ 経路制限付き隠れマルコフモデルでモデル化
 (離散時間表現にしたかったのはこの為)

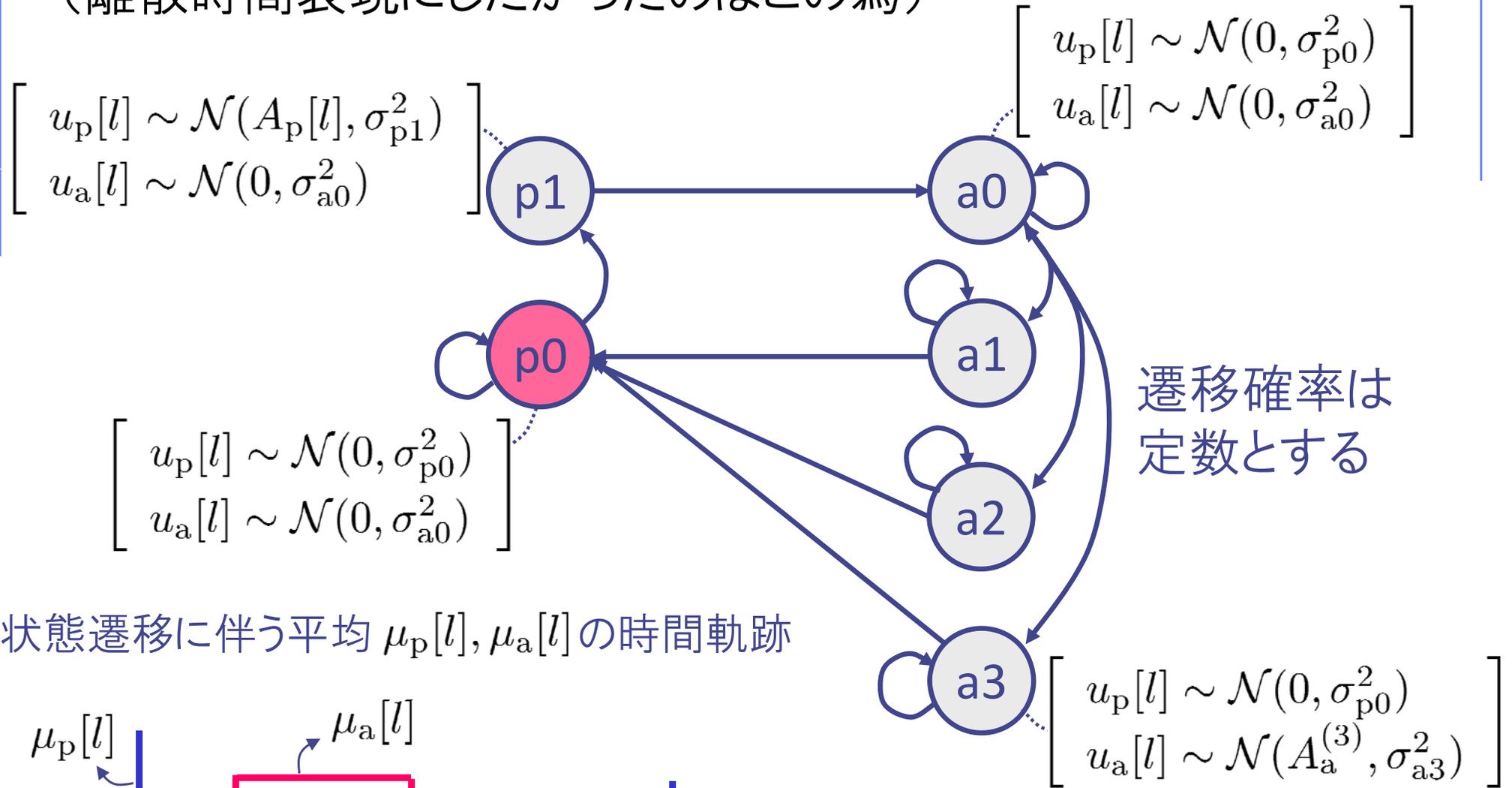


状態遷移に伴う平均 $\mu_p[l], \mu_a[l]$ の時間軌跡

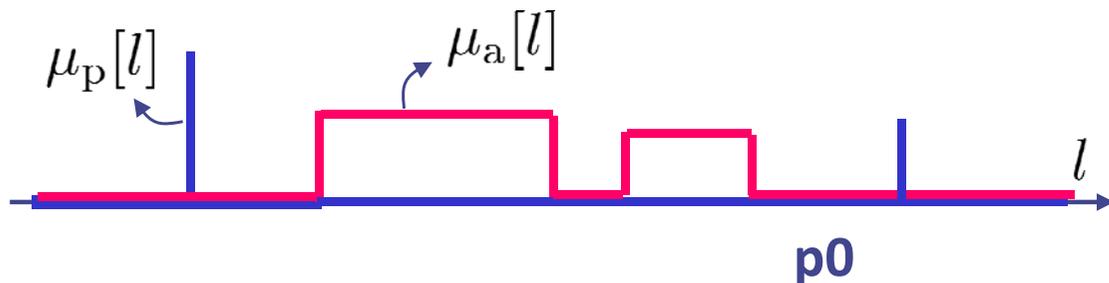


フレーズ指令 & アクセント指令の確率モデル化

◆ 経路制限付き隠れマルコフモデルでモデル化
 (離散時間表現にしたかったのはこの為)

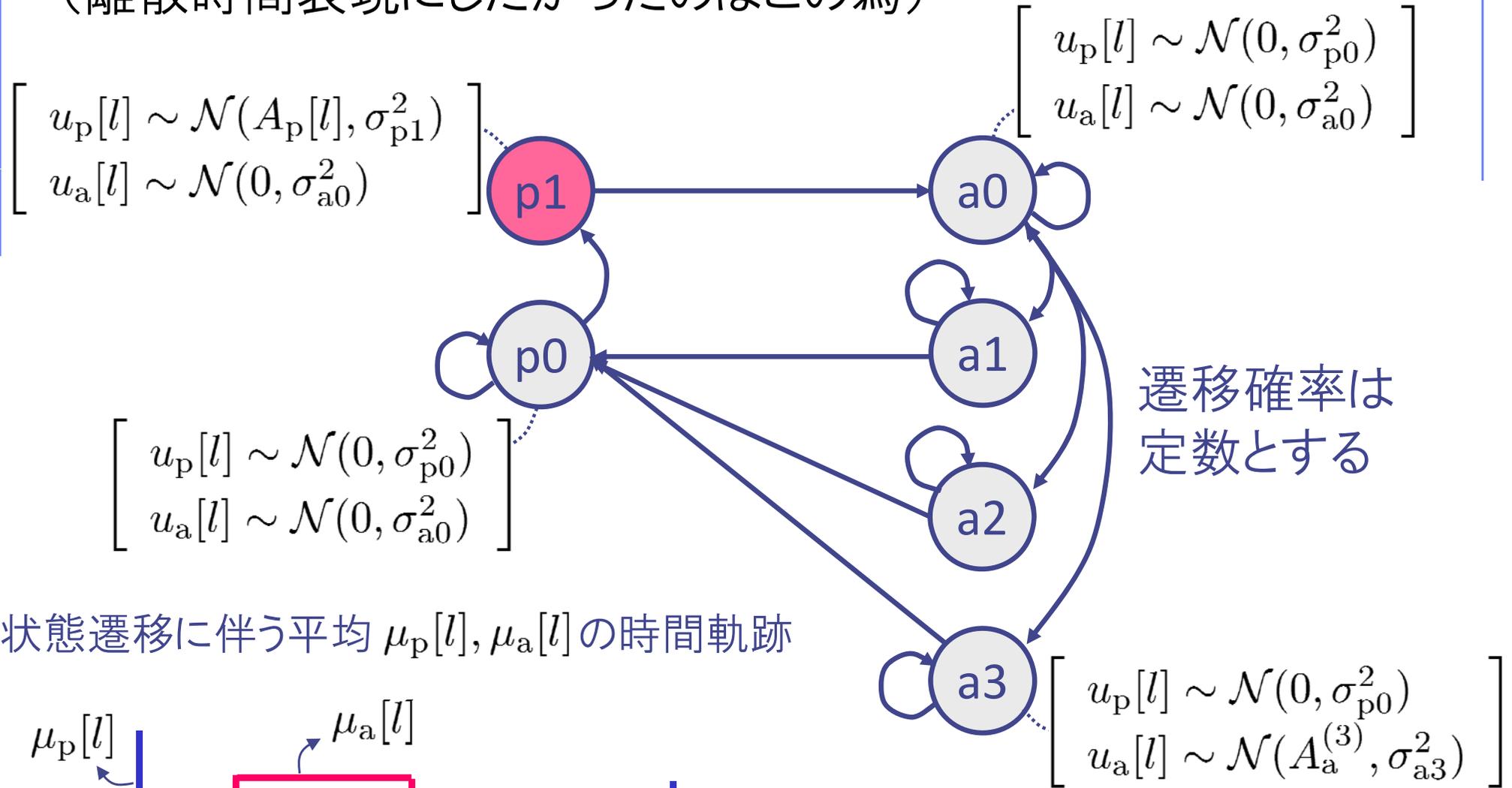


状態遷移に伴う平均 $\mu_p[l], \mu_a[l]$ の時間軌跡

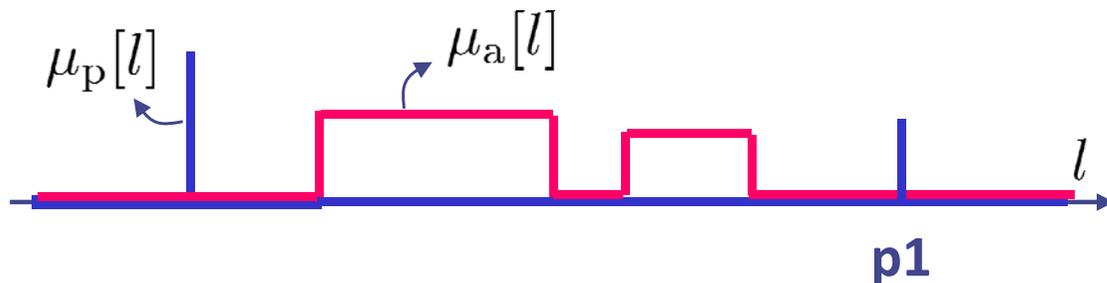


フレーズ指令 & アクセント指令の確率モデル化

◆ 経路制限付き隠れマルコフモデルでモデル化
 (離散時間表現にしたかったのはこの為)

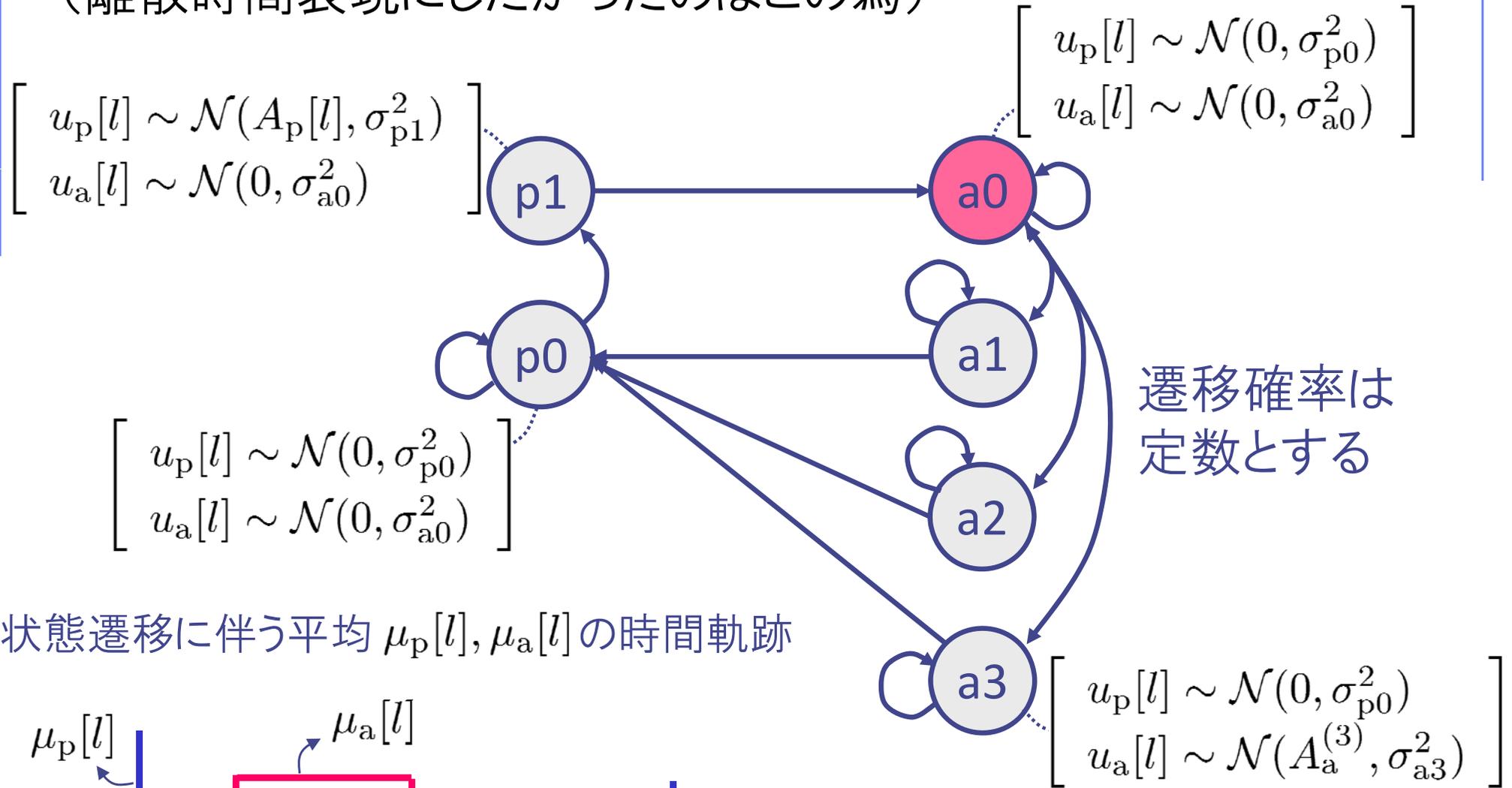


状態遷移に伴う平均 $\mu_p[l], \mu_a[l]$ の時間軌跡

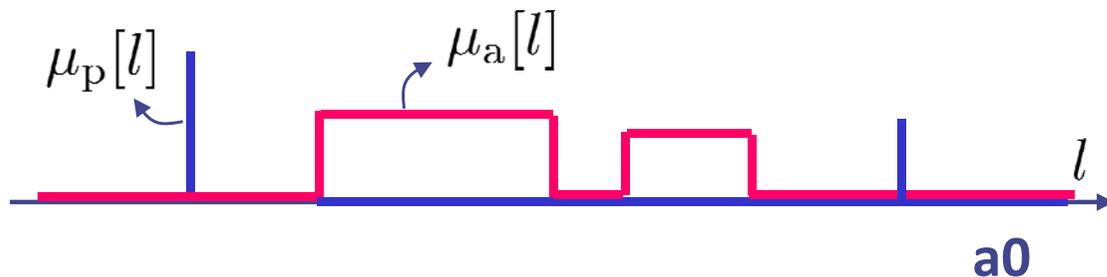


フレーズ指令 & アクセント指令の確率モデル化

◆ 経路制限付き隠れマルコフモデルでモデル化
 (離散時間表現にしたかったのはこの為)



状態遷移に伴う平均 $\mu_p[l], \mu_a[l]$ の時間軌跡



本研究の目的

(1) 藤崎モデルの離散時間表現を統計モデル化 (2)

- 統計的手法を駆使した強力なパラメータ推定法の枠組を確立
- 統計モデルに基づく音声処理問題(合成、分析、分離、強調)にスムーズに組み込める基盤を構築

発表アウトライン

◆モデル化の手順

(1) 後退差分近似による連続時間藤崎モデルの離散化

(2) フレーズ指令 & アクセント指令の確率モデル化

→ (1)と(2)からF0パターンの確率密度関数を導出

◆提案モデルの使用方法

◆パラメータ推定アルゴリズムの導出

本研究の目的

(1) 藤崎モデルの離散時間表現を統計モデル化

- 統計的手法を駆使した強力なパラメータ推定法の枠組を確立
- 統計モデルに基づく音声処理問題(合成、分析、分離、強調)にスムーズに組み込める基盤を構築

発表アウトライン

◆モデル化の手順

(1) 後退差分近似による連続時間藤崎モデルの離散化

(2) フレーズ指令 & アクセント指令の確率モデル化

→ (1)と(2)からF0パターンの確率密度関数を導出

◆提案モデルの使用方法

◆パラメータ推定アルゴリズムの導出

F0パターンの確率密度関数の導出

◆ フレーズ成分 $\Omega_p = (\Omega_p[1], \dots, \Omega_p[T])^T$ の確率密度関数

(1) で得た拘束式 $u_p[l] = g_0^p \Omega_p[l] + g_1^p \Omega_p[l-1] + g_2^p \Omega_p[l-2]$

(2) で得た確率モデル $u_p[l] \sim \mathcal{N}(\mu_p[l], \sigma_p[l]^2)$ 状態系列によって決まる
平均と分散の遷移系列

$$\begin{array}{c} \underbrace{\hspace{10em}}_{\mathbf{G}_p} \end{array}
 \begin{bmatrix} g_0^p & & & & 0 \\ g_1^p & g_0^p & & & \\ g_2^p & g_1^p & g_0^p & & \\ & \ddots & \ddots & \ddots & \\ 0 & & g_2^p & g_1^p & g_0^p \end{bmatrix}
 \begin{array}{c} \underbrace{\hspace{2em}}_{\Omega_p} \\ \begin{bmatrix} \Omega_p[1] \\ \Omega_p[2] \\ \Omega_p[3] \\ \vdots \\ \Omega_p[L] \end{bmatrix} \end{array}
 =
 \begin{array}{c} \underbrace{\hspace{2em}}_{u_p} \\ \begin{bmatrix} u_p[1] \\ u_p[2] \\ u_p[3] \\ \vdots \\ u_p[T] \end{bmatrix} \end{array}
 \sim \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$$

対角成分が $\sigma_p[l]^2$
の対角行列

➡ $\Omega_p \sim \mathcal{N}(\mathbf{G}_p^{-1} \boldsymbol{\mu}_p, \mathbf{G}_p^{-1} \boldsymbol{\Sigma}_p (\mathbf{G}_p^{-1})^T)$

◆ アクセント成分も同様 $\Omega_a \sim \mathcal{N}(\mathbf{G}_a^{-1} \boldsymbol{\mu}_a, \mathbf{G}_a^{-1} \boldsymbol{\Sigma}_a (\mathbf{G}_a^{-1})^T)$

◆ F0パターンベクトル $\Omega = \Omega_p + \Omega_a + \Omega_b$

➡ $\Omega \sim \mathcal{N}(\mathbf{G}_p^{-1} \boldsymbol{\mu}_p + \mathbf{G}_a^{-1} \boldsymbol{\mu}_a + \mu_b \mathbf{1},$

$\mathbf{G}_p^{-1} \boldsymbol{\Sigma}_p (\mathbf{G}_p^{-1})^T + \mathbf{G}_a^{-1} \boldsymbol{\Sigma}_a (\mathbf{G}_a^{-1})^T + \boldsymbol{\Sigma}_b)$

F0パターンの確率密度関数の導出

◆ フレーズ成分 $\Omega_p = (\Omega_p[1], \dots, \Omega_p[T])^T$ の確率密度関数

(1) で得た拘束式 $u_p[l] = g_0^P \Omega_p[l] + g_1^P \Omega_p[l-1] + g_2^P \Omega_p[l-2]$

(2) で得た確率モデル $u_p[l] \sim \mathcal{N}(\mu_p[l], \sigma_p[l]^2)$ 状態系列によって決まる
平均と分散の遷移系列

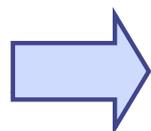
$$G_p = \begin{bmatrix} g_0^P & & & & \\ g_1^P & g_0^P & & & \\ g_2^P & g_1^P & g_0^P & & \\ & \dots & \dots & & \\ 0 & & g_2^P & & \end{bmatrix}$$

パラメータ Θ :

- 状態遷移系列 $\{s_l\}_{l=1}^L$
- 状態出力分布の平均値 $\{A_p[l]\}_{l=1}^L, \{A_a^{(n)}\}_{n=1}^N$
- フレーズ制御パラメータ ψ
- アクセント制御パラメータ φ
- ベースライン値 μ_b

◆ アクセント成分

◆ F0パターンベクトル $\Omega = \Omega_p + \Omega_a + \Omega_b$



$$\Omega \sim \mathcal{N}(G_p^{-1} \mu_p + G_a^{-1} \mu_a + \mu_b \mathbf{1},$$

$$G_p^{-1} \Sigma_p (G_p^{-1})^T + G_a^{-1} \Sigma_a (G_a^{-1})^T + \Sigma_b)$$

本研究の目的

(1) 藤崎モデルの離散時間表現を統計モデル化

- 統計的手法を駆使した強力なパラメータ推定法の枠組を確立
- 統計モデルに基づく音声処理問題(合成、分析、分離、強調)にスムーズに組み込める基盤を構築

発表アウトライン

◆モデル化の手順

(1) 後退差分近似による連続時間藤崎モデルの離散化

(2) フレーズ指令 & アクセント指令の確率モデル化

→ (1)と(2)からF0パターンの確率密度関数を導出

◆提案モデルの使用方法

◆パラメータ推定アルゴリズムの導出

本研究の目的

藤崎モデルの離散時間表現を統計モデル化

- 統計的手法を駆使した強力なパラメータ推定法の枠組を確立
- 統計モデルに基づく音声処理問題(合成、分析、分離、強調)にスムーズに組み込める基盤を構築

発表アウトライン

◆モデル化の手順

(1) 後退差分近似による連続時間藤崎モデルの離散化

(2) フレーズ指令 & アクセント指令の確率モデル化

→ (1)と(2)からF0パターンの確率密度関数を導出

◆提案モデルの使用方法

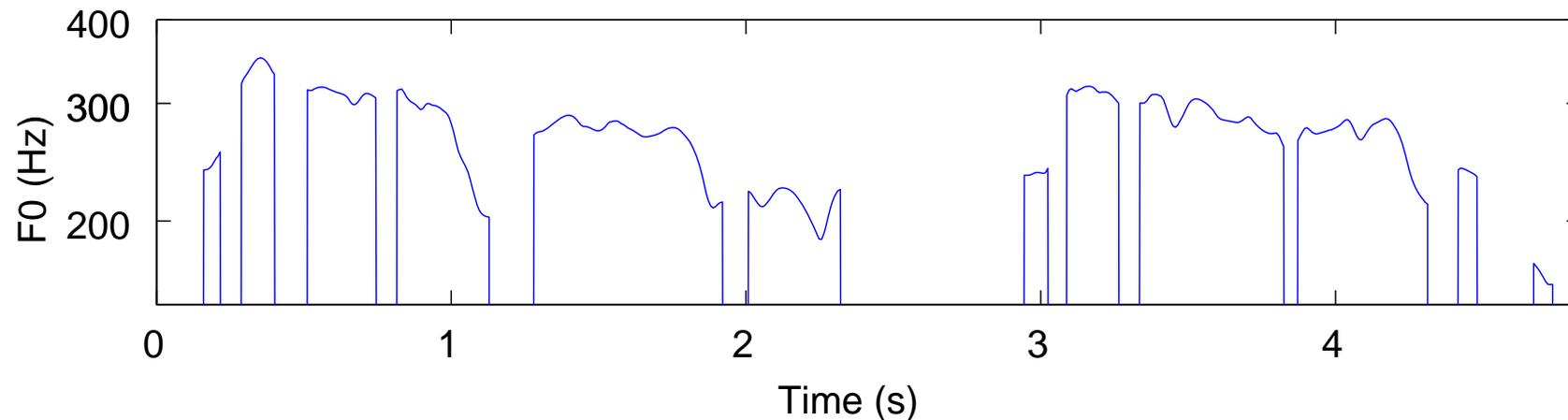
◆パラメータ推定アルゴリズムの導出

提案モデルの使用手法

$$\Omega \sim \mathcal{N}(\mathbf{G}_p^{-1} \boldsymbol{\mu}_p + \mathbf{G}_a^{-1} \boldsymbol{\mu}_a + \mu_b \mathbf{1}, \\ \mathbf{G}_p^{-1} \boldsymbol{\Sigma}_p (\mathbf{G}_p^{-1})^T + \mathbf{G}_a^{-1} \boldsymbol{\Sigma}_a (\mathbf{G}_a^{-1})^T + \boldsymbol{\Sigma}_b)$$

◆ 実測F0パターンからの藤崎モデルの推定

■ 実測F0パターンは”欠損データ”



- 上記統計モデルは全区間のデータが観測されていることが想定
→ 欠損データの下での最尤/最大事後確率推定は
EMアルゴリズムで解ける！

◆ F0をパラメータにもつ音声信号/スペクトルの統計モデルにおけるF0パラメータの事前分布として利用可能(詳細は直後の発表にて)

本研究の目的

藤崎モデルの離散時間表現を統計モデル化

- 統計的手法を駆使した強力なパラメータ推定法の枠組を確立
- 統計モデルに基づく音声処理問題(合成、分析、分離、強調)にスムーズに組み込める基盤を構築

発表アウトライン

◆モデル化の手順

(1) 後退差分近似による連続時間藤崎モデルの離散化

(2) フレーズ指令 & アクセント指令の確率モデル化

→ (1)と(2)からF0パターンの確率密度関数を導出

◆提案モデルの使用方法

◆パラメータ推定アルゴリズムの導出

本研究の目的

藤崎モデルの離散時間表現を統計モデル化

- 統計的手法を駆使した強力なパラメータ推定法の枠組を確立
- 統計モデルに基づく音声処理問題(合成、分析、分離、強調)にスムーズに組み込める基盤を構築

発表アウトライン

◆モデル化の手順

(1) 後退差分近似による連続時間藤崎モデルの離散化

(2) フレーズ指令 & アクセント指令の確率モデル化

→ (1)と(2)からF0パターンの確率密度関数を導出

◆提案モデルの使用方法

◆パラメータ推定アルゴリズムの導出

パラメータの最大事後確率推定

◆ Ω が与えられた下で $P(\Theta|\Omega) \propto \underline{P(\Omega|\Theta)}P(\Theta)$ を最大化

$$\Omega \sim \mathcal{N}(\mathbf{G}_p^{-1}\boldsymbol{\mu}_p + \mathbf{G}_a^{-1}\boldsymbol{\mu}_a + \mu_b\mathbf{1}, \mathbf{G}_p^{-1}\boldsymbol{\Sigma}_p(\mathbf{G}_p^{-1})^T + \mathbf{G}_a^{-1}\boldsymbol{\Sigma}_a(\mathbf{G}_a^{-1})^T + \boldsymbol{\Sigma}_b)$$

◆ フレーズ成分 Ω_p , アクセント成分 Ω_a , ベースライン成分 Ω_b を完全データと見なすとEM法が適用できる!

$$\mathbf{x} = \begin{bmatrix} \Omega_p \\ \Omega_a \\ \Omega_b \end{bmatrix} \sim \mathcal{N} \left(\underbrace{\begin{bmatrix} \mathbf{G}_p^{-1}\boldsymbol{\mu}_p \\ \mathbf{G}_a^{-1}\boldsymbol{\mu}_a \\ \mu_b\mathbf{1} \end{bmatrix}}_{\mathbf{m}}, \underbrace{\begin{bmatrix} \mathbf{G}_p^{-1}\boldsymbol{\Sigma}_p(\mathbf{G}_p^{-1})^T & O & O \\ O & \mathbf{G}_a^{-1}\boldsymbol{\Sigma}_a(\mathbf{G}_a^{-1})^T & O \\ O & O & \boldsymbol{\Sigma}_b \end{bmatrix}}_{\boldsymbol{\Lambda}} \right)$$

■ “不完全データ” Ω と完全データの関係: $\Omega = \underbrace{[I \ I \ I]}_H \begin{bmatrix} \Omega_p \\ \Omega_a \\ \Omega_b \end{bmatrix}$

■ Q関数

$$Q(\Theta, \Theta') \stackrel{c}{=} -\frac{1}{2} \left[\text{tr}(\boldsymbol{\Lambda}^{-1} \mathbb{E}[\mathbf{x}\mathbf{x}^T | \Omega; \Theta]) - 2\mathbf{m}^T \boldsymbol{\Lambda}^{-1} \mathbb{E}[\mathbf{x} | \Omega; \Theta] + \mathbf{m}^T \boldsymbol{\Lambda}^{-1} \mathbf{m} \right]$$

EMアルゴリズム

◆ E-step

F0パターンモデル

$$\begin{aligned}\mathbb{E}[x|\Omega; \Theta] &= m + \Lambda H^T (H \Lambda H^T)^{-1} (\Omega - \underline{Hm}) \\ \mathbb{E}[xx^T|\Omega; \Theta] &= \Lambda - \Lambda H^T (H \Lambda H^T)^{-1} H \Lambda + \mathbb{E}[x|\Omega; \Theta] \mathbb{E}[x|\Omega; \Theta]^T\end{aligned}$$

◆ M-step

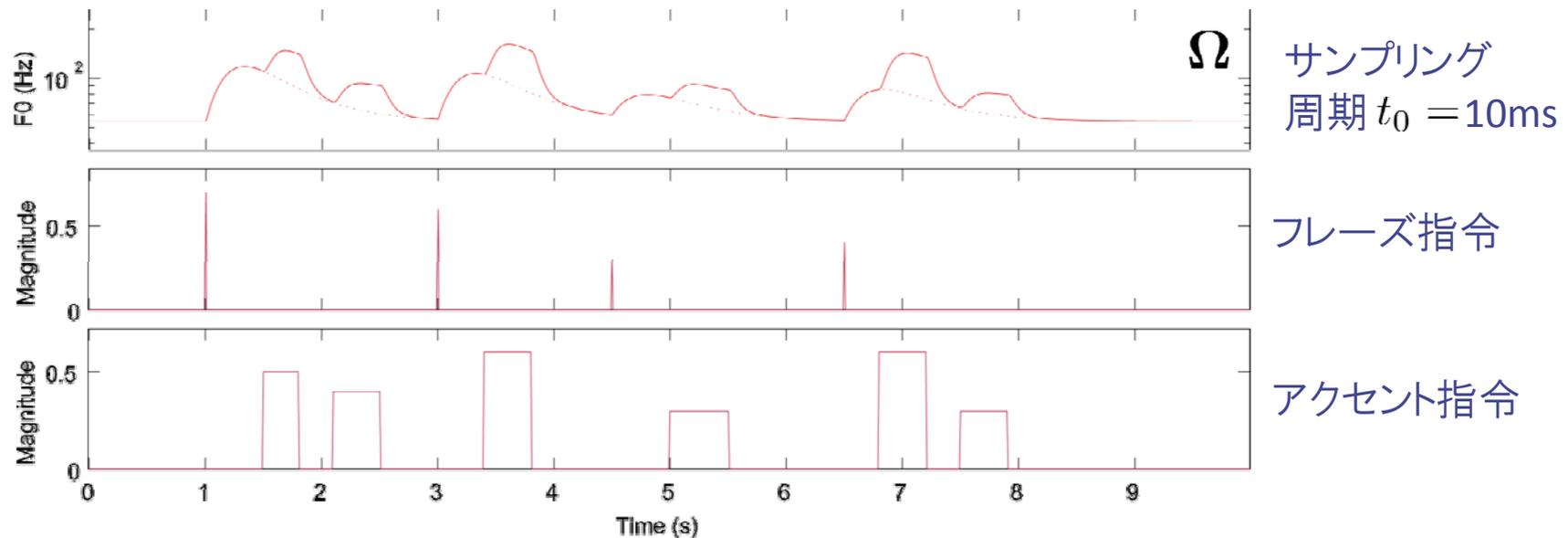
$$Q(\Theta, \Theta') \stackrel{c}{=} -\frac{1}{2} \left[\text{tr}(\Lambda^{-1} \mathbb{E}[xx^T|\Omega; \Theta]) - 2m^T \Lambda^{-1} \mathbb{E}[x|\Omega; \Theta] + m^T \Lambda^{-1} m \right]$$

を Θ に関して最大化

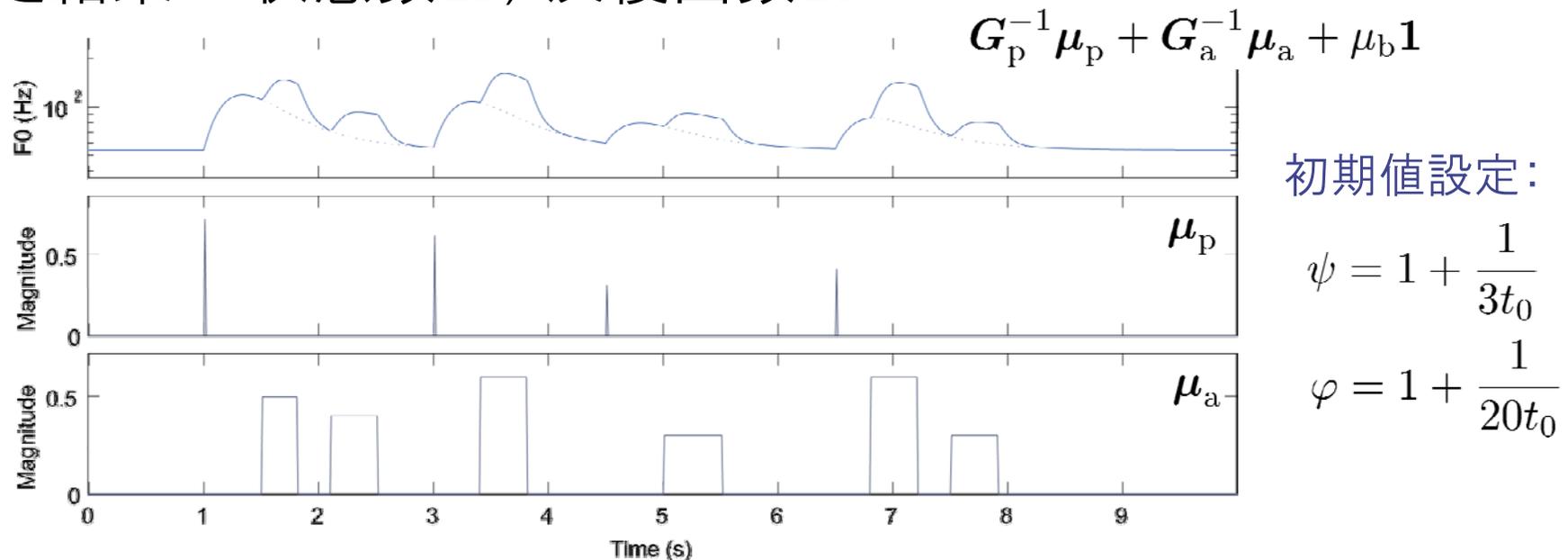
- ψ, φ の更新は4次方程式の求解に帰着
- 状態遷移系列と状態出力分布の平均値の更新はViterbi学習により実現可能

シミュレーション実験

◆ 実験データ: 元の藤崎モデルから合成した人工データ



◆ 推定結果: 状態数13, 反復回数10



まとめ

◆ 藤崎モデルの離散時間表現の確率モデルを提案

- 統計的手法を駆使した強力なパラメータ推定法の枠組を構築
- 音声の統計モデルに基づく各種音声処理問題にスムーズに組み込める基盤を構築

◆ 提案モデルのポイント

- もともとは連続時間系の藤崎モデルを離散時間系で記述
- フレーズ指令とアクセント指令を隠れマルコフモデルにより確率モデル化
- EMアルゴリズムによるパラメータ推定法を導出

◆ 効果

- 今のところ人工データに対する一有効性だけ確認済
- 実測F0パターンに対してはどうか？ 欠損データの影響は？
- 統計的音声モデルに組み込んだ際の効果は？

◆ 参考までに

- 3-P-31 複数振動基底に基づく歌声のF0 動特性の統計的モデリング
©大石 康智, 亀岡 弘和, 持橋 大地, 永野 秀尚, 柏野 邦夫