

物語イベントの状態遷移解析：半教師あり GMM を用いた 転調点抽出

——宮沢賢治『銀河鉄道の夜』を例に

要旨

本研究は、物語を「状態」と「遷移」として定式化し、転調点を再現可能な量として抽出する枠組みを提示する。宮沢賢治『銀河鉄道の夜』のイベント列を、心理特徴量 $X=(m, iso)$ へ写像し、少量ラベル（アンカー）を用いた半教師あり GMM を EM で推定する。各イベントの事後確率エントロピーを boundaryness と定義し、上位点を転調点候補として抽出する。モデル選択（BIC/AIC/Silhouette）と alpha 比較（アンカー整合、境界安定、遷移数、entropy 総量）により、二状態モデルが解釈性と指標の両面で最適であることを示す。さらに、転調点は遷移回数と一致せず、心理的不確定性の極大として分離して現れることを明らかにする。

キーワード：半教師あり学習、ガウス混合モデル、EM アルゴリズム、転調点、エントロピー、デジタル・ヒューマニティーズ

1. 背景

物語を「状態」と「遷移」として定式化する試みは、解釈の再現可能性を担保するデジタル・ヒューマニティーズの基盤である。しかし現状の物語解析は、(i) 章・場面の区切りを人手に依存する、(ii) 分割点の意味づけが事後的で恣意的になりやすい、(iii) 教師データ不足により監督学習が成立しにくい、という制約を併せ持つ。この制約は、局面の転調が重要である一方でラベル付けが高コストであり、かつ解釈可能性が要求される文学作品で顕在化する。ゆえに、少量のラベルで状態の意味を固定しつつ、転調点をモデル内部量として抽出できる手法が必要である。

1.1 関連研究（4つの束）

本研究の位置づけを明確化するため、先行研究を（A）デジタル・ヒューマニティーズ／計算ナラトロジー、（B）物語可視化・プロット曲線、（C）物語を低次元軌跡として捉える時系列モデル、（D）半教師あり／アンカークラスタリングと不確実性指標、の4つに整理する。

（A）デジタル・ヒューマニティーズ／計算ナラトロジー：計算モデルは、物語生成や理解を「研究のための作動する表現」として実装し、理論検証と洞察獲得に用いられる（Montfort & Pérez y Pérez, 2023）。また CMN（Computational Models of Narrative）系の研究コミュニティでは、イベント単位の表現、注釈付きコーパス、再現性の確保、プロット曲線の利用などが継続的に議論されてきた（CMN'12 Proceedings）。一方で、物語イベントの切り分けは語用・修辞・読者経験に強く依存し、離散化は人工的構造を持ち込みやすいという批判も整理されている（DHQ 系の総説：000548）。したがって、（i）イベント列を明示的に定義し根拠引用を保持すること、（ii）解釈の主観性を測定可能な形で管理することが、方法論上の要件となる。

（B）物語可視化・教育利用：PlotVis や CATMA、Narrelations など、注釈に基づき物語構造を可視化して学習・比較に用いる試みが報告されている（000548）。これらは「解釈の多様性を許容しつつ、議論可能な形式に落とす」点で有用だが、スケール（イベント粒度）とスケーラ

ビリティ（多作品への拡張）のトレードオフが残る（000548）。本研究の“章×イベント列×2軸”の設計は、このトレードオフに対し、最小限の注釈で比較可能な可視化を得ることを狙う。

（C）低次元軌跡としての物語：物語を感情・価値・緊張などの連続量の軌跡として表現し、典型形状（アーク）や転換点を議論する研究がある（例：Reagan et al., 2016）。この流れは、文学解釈を“時間的変化の形”として比較可能にする。さらに、物語生成・理解における「感情テンションの上昇下降」など、曲線形状を内部評価に使う発想も提示されている（CMN'12 Proceedings）。本研究は、m/isoの2軸を“心理的状态空間”として置き、GMMで潜在状態を推定し、曲線ではなく「状態遷移＋不確実性スパイク」として転調を定義する点で、この系譜を拡張する。

（D）半教師あり／アンカークラスタリングと不確実性：少数のラベル（アンカー）でクラスタの意味を固定し、未ラベル点をEMで推定する枠組みは、解釈可能性とデータ不足の両立を狙う実践的アプローチである。機械学習側では、アンカー（信頼できる代表点）を用いてクラスタ中心を導き、擬似ラベルの品質をKL距離などでフィルタする「anchored clustering」系の手法が提案されている（NeurIPS 2022）。本研究の α 重み付けは、この“アンカーで意味を固定しつつ、残りは確率的に委ねる”設計思想に対応する。また境界点（boundaryness）を事後確率エントロピーで測るのは、割り当てが拮抗する点を機械的に抽出し、読解の焦点として返す仕掛けである。

以上より、本研究はAの「計算モデルを人文学的探究に使う」立場を取りつつ、B/Cの可視化・軌跡研究に接続し、Dの半教師ありクラスタリングと不確実性指標を“解釈の再現可能性”のために導入する点に特徴がある。

1.2 本究の貢と新規性

本研究の貢献は次の3点である。(1) 根拠引用付きイベント列（JSON）を中心に、解釈の単位と証拠を固定することで、物語分析の再現性問題を操作可能にした。(2) m/isoの2軸を用いた最小表現により、章をまたぐ比較と可視化を可能にした。(3) 半教師ありGMMにアンカーを導入し、さらに事後確率エントロピーを転調点指標として用いることで、「どこが揺れているか」を機械的に抽出し、本文への還元（evidence参照）を可能にした。

2. 手法

2.1 データ（イベント列）

入力章（scene_id）ごとに分割されたイベント列であり、各イベントに根拠引用（evidence）と要約（desc）を付与する。イベントは全体時系列（global_step）として連番管理し、章開始位置は時系列上の参照線として利用する。

2.1.1 アノテーション手順（運用プロンプトの要約）

本研究では、各章本文のみを根拠としてイベント列を作成した。運用上は、以下の簡易ルールをプロンプトとして固定した。

- 1 event は「ジョバンニの状態が更新される」まとめり（目安 1〜4 文）。章あたり 6〜30 events を基本とし、文章量が多い章（例：第 9 章）は粒度を細かくして event 数を増やした。
- event 名は中立にし、原文にない情報は追加しない（推測が必要な場合は desc に「推測」と明記）。
- 各 event に evidence（原文短引用）を 1〜2 本付与（各引用は 40 文字以内、改変しない）。
- 各 event に 2 軸特徴 $x=[\text{Morality, Isolation/Survival Urgency, local_time}]$ を付与（0〜1）。
- global_step は start_global_step から 1 ずつ増加。local_step も保持し、 $\text{local_time}=(\text{local_step}-1)/(\text{local_step_max}-1)$ で正規化。全章統合後に global_time を再計算できるようにする。
- アンカー（labeled）は「解釈の余地が小さい」イベントのみに限定し、章あたり最大 2 件（Communal_Happiness / Self_Protection）。迷う場合は unlabeled とする。

2.1.2 アンカー選定と簡易チェック（運用ルール）

本研究はアノテーションを単一の AI（固定プロンプト）で実行し、各イベントに原文引用（evidence）を必須化することで、後から第三者が照合できる監査可能性（audit trail）を確保した。実務上は、アンカーと転調点候補のみを人間がスポットチェックすることで、コストを抑えつつ品質を担保できる。

- アンカー（labeled）は「解釈の余地が小さい」点のみに限定（章あたり最大 2 件）。
- 選定基準：①ループリック上で極端（m が高い/低い、iso が高い等）② evidence が明確で短引用で示せる③章内で偏らない（可能な範囲で分散）。
- 簡易チェック：人間は（i）アンカーの妥当性（ラベルの明白さ）と（ii）entropy ピーク周辺の evidence 妥当性のみを確認し、必要なら小さな修正に留める。
- 将来的には、別モデル（例：別 LLM）に同一プロンプトで注釈させ、Top-K 転調点やアンカー一致度を比較することで、モデル間再現性を評価できる。

補足（本ケースのスポットチェック結果）：著者が抽出済みのアンカー一覧を本文（evidence）と照合したところ、全アンカーが作品の趣旨・価値観・設定と整合し、修正を要しなかった。この結果は定量的一致度ではないが、アンカーの顔妥当性（face validity）を支持する。

2.2 特徴量設計

各イベントを、以下の 2 次元心理特徴量で表現する。

- m（Morality）：自己防衛・自己中心への退避（0）から、共感・他者配慮・共同幸福（1）への連続値。
- iso（Isolation/Survival Urgency）：安心・切迫感の低さ（0）から、孤独・欠乏・喪失・切迫（1）への連続値。

- Time（参照軸）：global_step を正規化して 0-1 に写像し、可視化の参照軸として用いる。

2.3 半教師あり GMM（EM）

観測 $X=(m, iso)$ に対し、 K 成分の GMM を EM で推定する。少数の labeled 点（アンカー）を与え、unlabeled 点と混合して推定する。アンカーの寄与は係数 α で重み付けし、 α が大きいほどクラスタ意味（解釈ラベル）が固定化される。推定は複数初期値で行い、seed に対する安定性を評価する。

2.3.1 初学者向け補足：GMM/EM（アンカー付き）と boundaryness の意味

本節では、2.3～2.4 で用いる GMM/EM と、アンカー（半教師あり）および boundaryness（entropy）を、初学者が再現できる粒度で要点整理する。

（1）GMM の生成モデル

各イベント時刻 t の特徴ベクトル $x_t \in R^d$ （本研究では $d=2$ ： m, iso ）を、 K 個の潜在状態 $z_t \in \{1, \dots, K\}$ から生成される混合ガウス分布で表す。混合比 ϕ_k （ $\sum_k \phi_k = 1$ ）、平均 μ_k 、共分散 Σ_k をパラメータ $\theta = \{\phi, \mu, \Sigma\}$ とする。

生成過程：

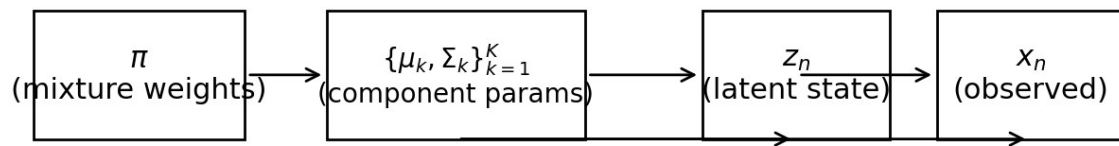
- 1) $z_t \sim \text{Categorical}(\phi)$
- 2) $x_t | (z_t=k) \sim N(\mu_k, \Sigma_k)$

周辺尤度（混合分布）：

$$p(x_t | \theta) = \sum_{k=1..K} \phi_k \cdot N(x_t | \mu_k, \Sigma_k)$$

目的は対数尤度の最大化：

$$L(\theta) = \sum_{t=1..T} \log \left(\sum_{k=1..K} \phi_k \cdot N(x_t | \mu_k, \Sigma_k) \right)$$



Generative story:

- 1) Draw $z_n \sim \text{Categorical}(\pi)$
- 2) Draw $x_n \sim N(\mu_{z_n}, \Sigma_{z_n})$

図 2.3.1-1 : GMM の生成モデル (模式図)

(2) EM アルゴリズム (E-step / M-step)

GMM は「log(和)」を含むため直接最適化が難しい。そこで潜在変数 z を導入し、EM で反復的に尤度を改善する。

E-step : 責務 (responsibility) の計算

$$w_{tk} = p(z_t = k | x_t, \theta) = \phi_k N(x_t | \mu_k, \Sigma_k) / \sum_j \phi_j N(x_t | \mu_j, \Sigma_j)$$

M-step : 重み付き十分統計量で更新

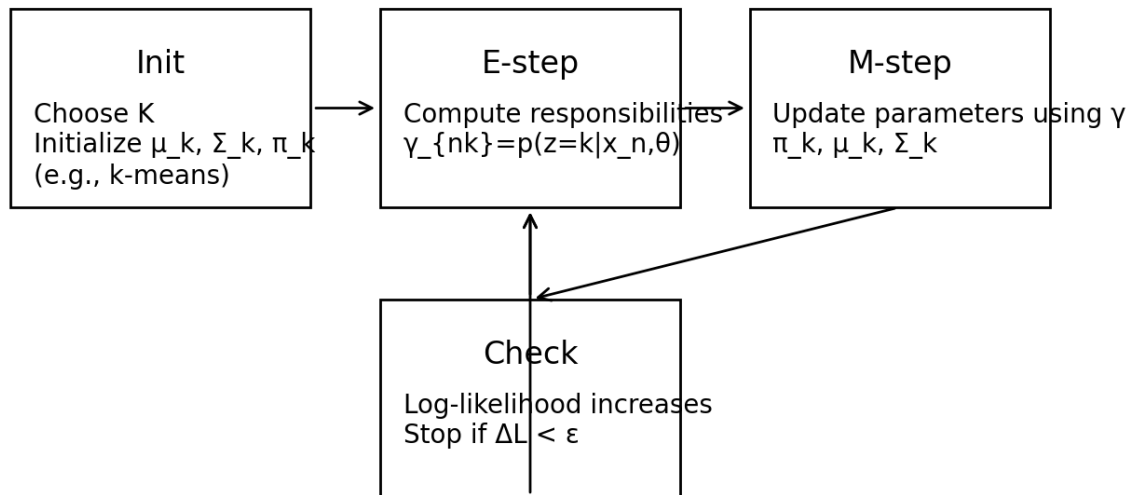
$$N_k = \sum_t w_{tk}$$

$$\phi_k = N_k / \sum_j N_j$$

$$\mu_k = (\sum_t w_{tk} x_t) / N_k$$

$$\Sigma_k = (\sum_t w_{tk} (x_t - \mu_k)(x_t - \mu_k)^T) / N_k$$

実装では数値安定化のため Σ_k の対角に正則化 reg を加え、 $\text{diag}=\text{True}$ の場合は対角共分散に制約する。



EM guarantees non-decreasing log-likelihood for GMM.

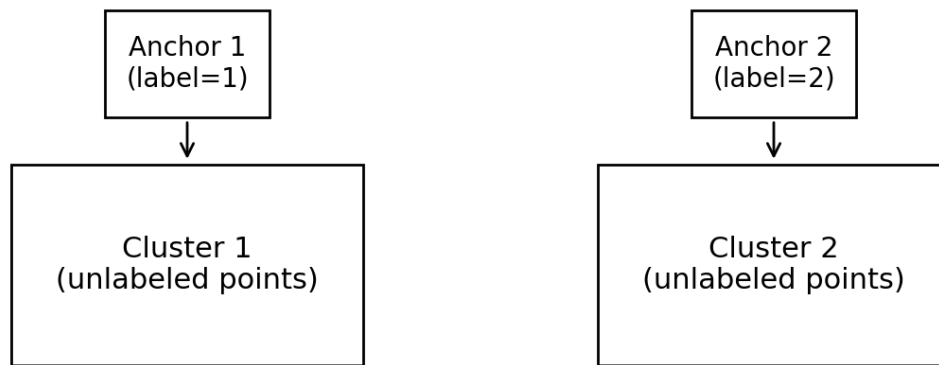
図 2.3.1-2 : EM の反復手順 (模式図)

(3) アンカー（半教師あり）と α の実装上の意味

本研究の半教師あり GMM は、未ラベル集合 X_u に対して E-step で w_u を計算し、M-step では X_u とラベル付き集合（アンカー） X_l を結合して更新する。アンカーは one-hot 行列 Y ($Y_{ak}=1[y_a=k]$) として表し、 $w_{all}=[w_u; \alpha Y]$ を用いて重み付き M-step を行う。すなわち、アンカー点は「既知ラベルへの割当を α 倍の擬似カウントとして注入」する方式であり、 α はクラスタの意味（解釈）を固定する強さを制御する。 $\alpha=0$ のとき教師なし GMM へ退化する。

実装に対応する更新（M-step の入力）：

$X_{all}=[X_u; X_l]$, $w_{all}=[w_u; \alpha Y]$ を作り、 $m_step_weighted(X_{all}, w_{all})$ で (ϕ, μ, Σ) を更新する。



Idea: treat anchor points as labeled observations with extra weight α .
 Add $\alpha \cdot \log p(x_a | z=y_a, \theta)$ to the objective.
 Cluster meaning is stabilized while other points remain soft-assigned.

図 2.3.1-3：アンカーと α の直感（模式図）

（4）boundaryness（entropy）を転調点とみなす根拠

学習後、各イベント i の事後確率 $w_i = p(z|x_i)$ が得られる。状態が明確な区間では w_i は一つのクラスタに集中し、状態が揺れる区間では複数クラスタに分散する。本研究ではこの「揺れ」を事後確率エントロピーで定量化し、boundaryness と定義する。

$$H_i = - \sum_{k=1..K} w_{ik} \log w_{ik}$$

H_i が大きい点は「どの状態にも属し切らない」曖昧区間であり、物語の葛藤・逡巡・価値衝突などが局所的に現れやすい。そこで H_i 上位（Top20 など）を転調点候補として抽出し、evidence（原文引用）へ還元して解釈する。

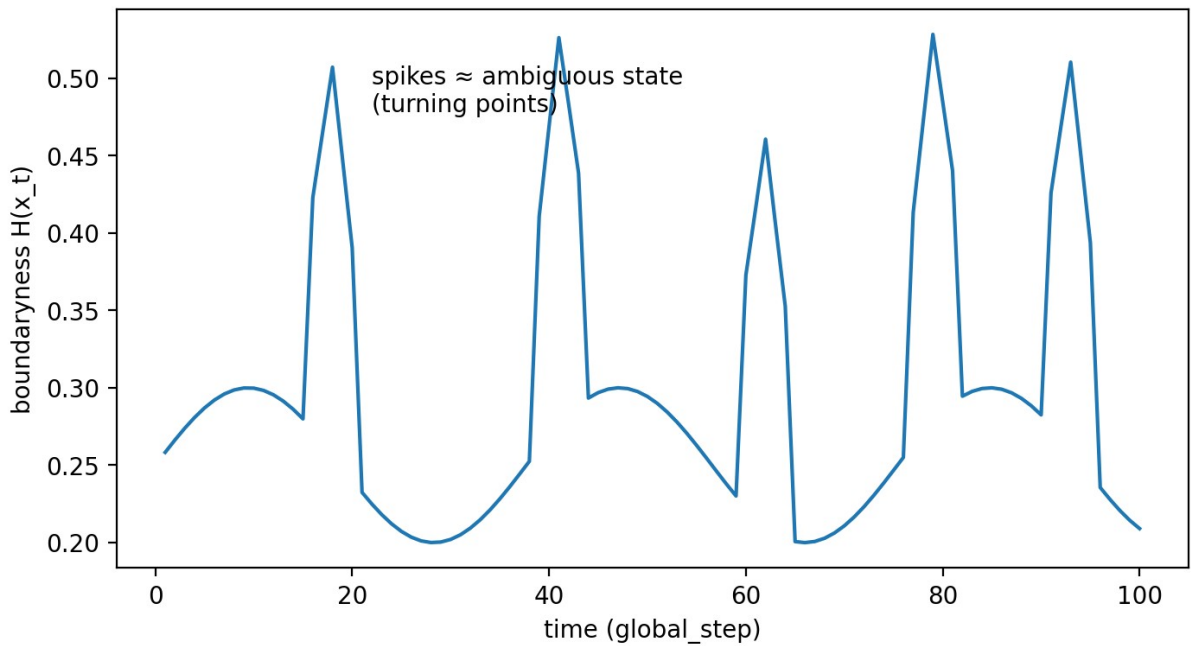


図 2.3.1-4 : entropy スパイクと転調点 (模式図)

2.4 転調点 (boundary) の定義

各イベント i について GMM の事後確率 $p(z|x_i)$ を算出し、そのエントロピー H_i を boundaryness と定義する。 H_i が大きい点は状態帰属が曖昧な点であり、転調点候補である。本稿では便宜上、entropy 上位 20 点（全体 98 点中、約 20%）を boundary として抽出し、時系列上に配置する。閾値（Top-K/上位割合）はハイパーパラメータであり、感度分析は今後の課題とする。

$$H_i = - \sum_{k=1..K} p(z=k | x_i) \log p(z=k | x_i)$$

2.5 評価指標

(i) K 選択：BIC/AIC/Silhouette により適合と分離を比較し、目的（転調抽出）に照らして解釈可能性を優先する。(ii) alpha 効果：アンカー整合、boundary の seed 安定性、alpha=0 に対する境界重なり、遷移回数、entropy 総量を比較する。

表 1: ルーブリック（m/iso/Time の定義と採点規則）

特徴量	0 の意味	1 の意味
m (Morality)	0=自己防衛・自分中心への退避	1=共感・他者救済・自己犠牲の受容
iso (Isolation/ Survival_Urgency)	0=安定（切迫感が低い）	1=孤独・欠乏・喪失が差し迫る（存在的危機）
Time（参考）	0=物語の冒頭	1=物語の終盤（global_step 正規化）

表 2: K 別モデル選択指標（BIC/AIC/Silhouette/Anchor agreement）

K	BIC (↓)	AIC (↓)	Silhouette (↑)	Anchor agreement
1.0	-197.66039	-208.00026	nan	0.714286
2.0	-197.834997	-221.099705	0.410386	0.857143
3.0	-182.828083	-219.017627	0.362225	1.0
4.0	-175.279388	-224.39377	0.245848	0.857143
5.0	-156.319921	-218.359141	0.362089	1.0
6.0	-140.541032	-215.505089	0.38262	1.0

表 3: alpha 別比較 (anchor_acc, boundary_Jaccard, transitions, entropy_sum 等)

alpha	anchor_acc	boundary_Jaccard	overlap_v s0	transitions	entropy_sum	seeds
0.0	0.857	1.0	1.0	14.0	4.244	5.0
1.0	1.0	0.76	0.495	20.4	16.665	5.0
5.0	1.0	1.0	0.905	12.0	3.694	5.0
10.0	1.0	1.0	0.905	12.0	3.152	5.0
20.0	1.0	1.0	0.818	12.0	3.199	5.0
50.0	1.0	1.0	0.667	16.0	3.47	5.0

図 1: m-iso 散布図 (点サイズ=entropy、円=boundary Top20)

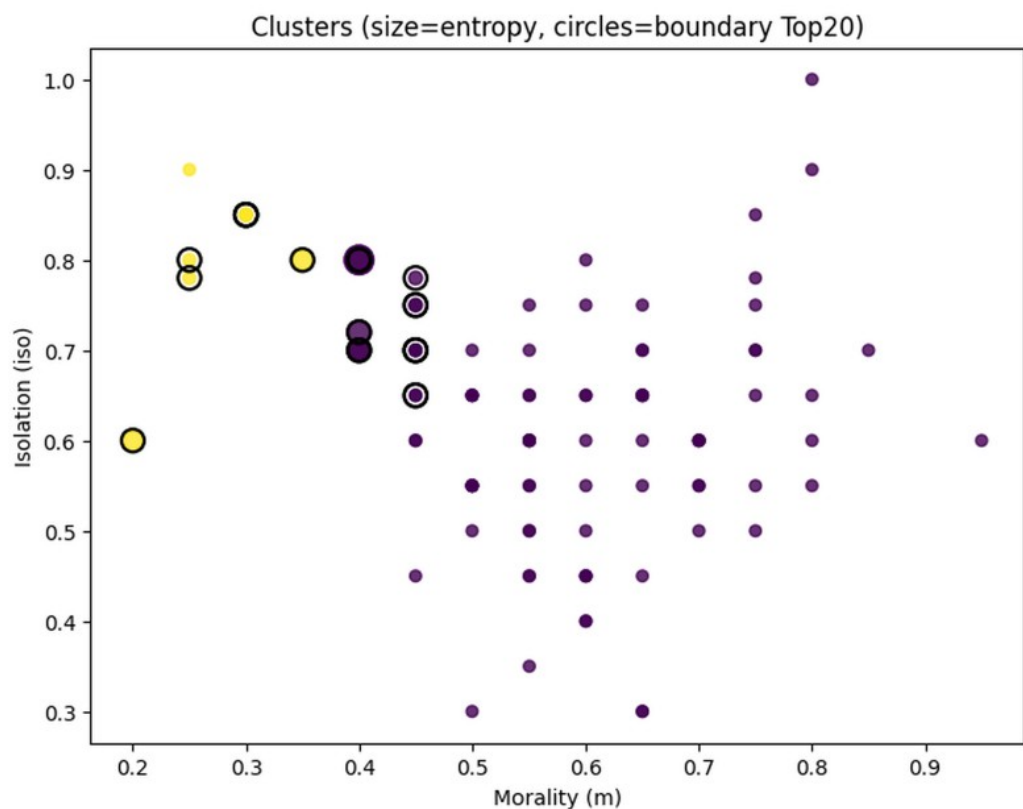


図 2: クラスタ中心と共分散楕円 (2σ)

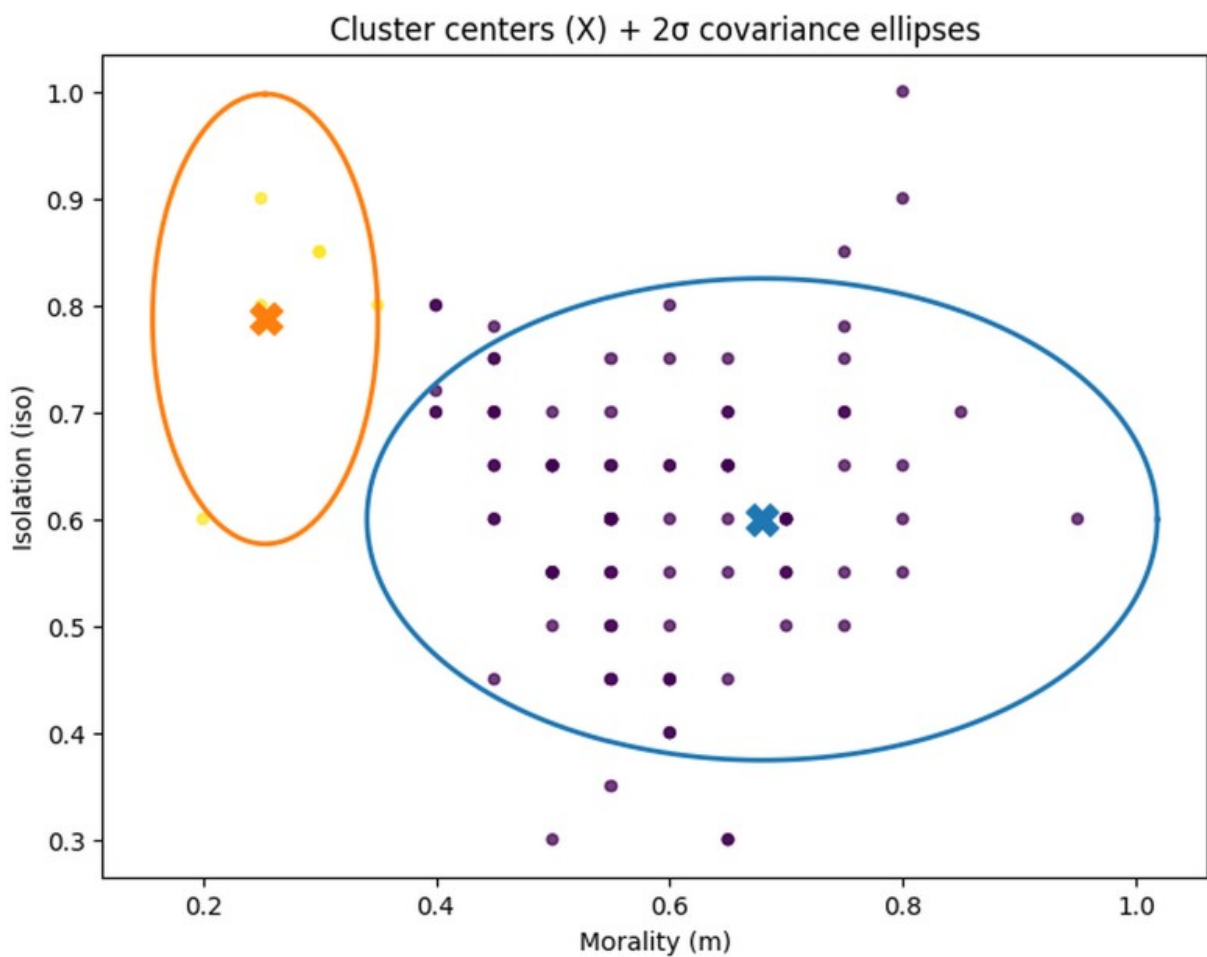


図 3: boundaryness 時系列 (縦線=章開始)

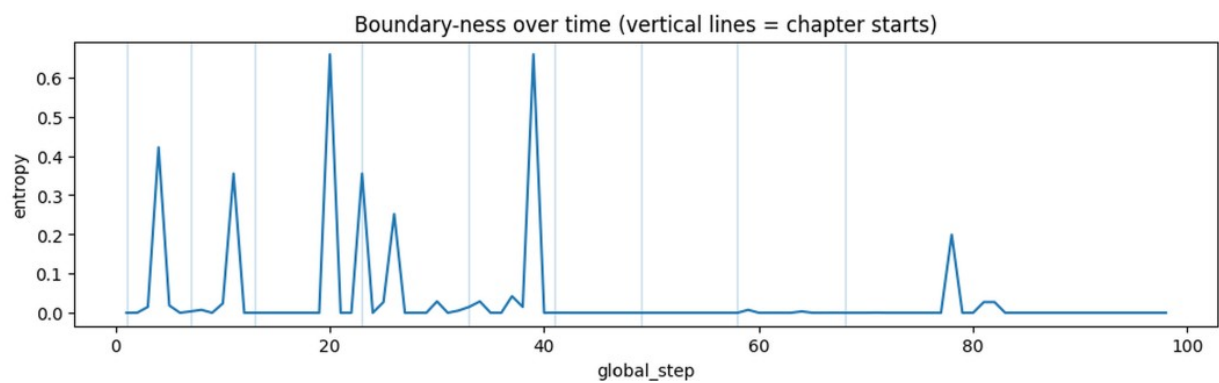


図 4: 状態 (クラスタ意味) の時系列

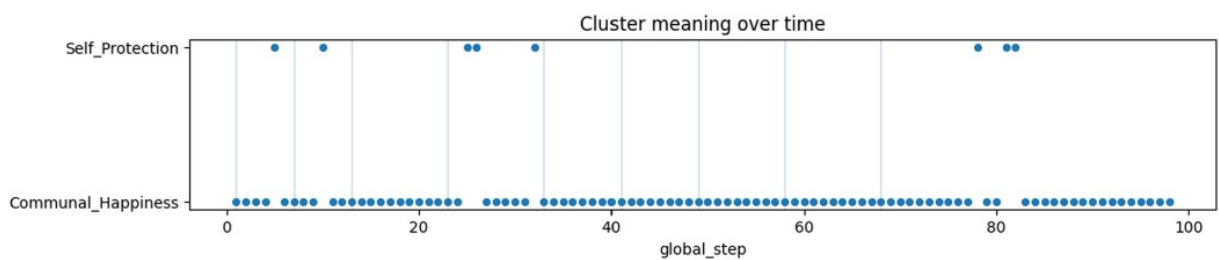


図 5: 章別 turning density (entropy 総量)

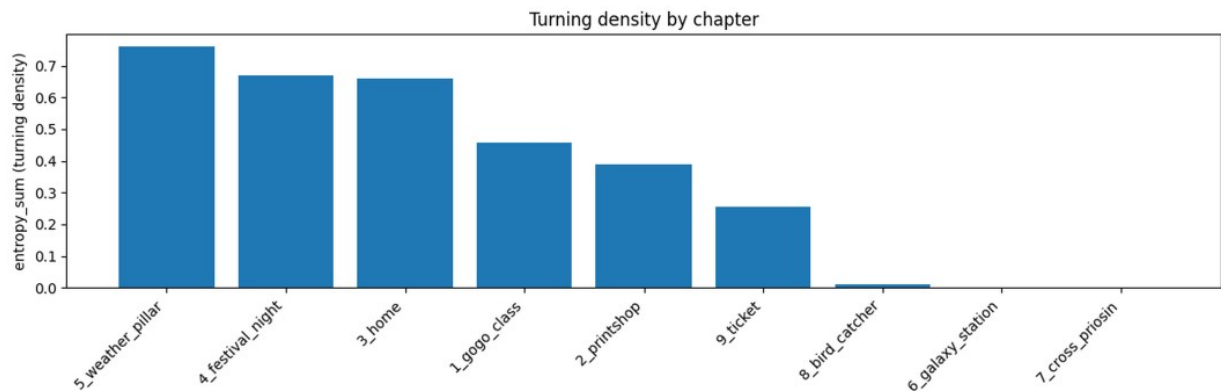
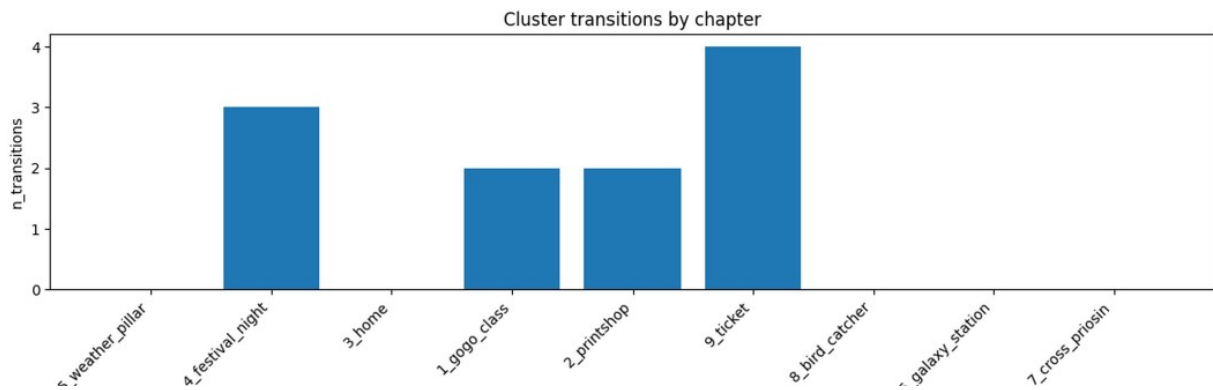


図 6: 章別クラスタ遷移回数



3. 結果と考察

3.1 本データ（2次元特）では $K=2$ が最も解可能である

$K=1$ は状態分割が生じず、遷移・転調の定義が成立しない。 $K \geq 3$ は 2 次元特徴 (m/iso) 上での過剰分割により、転調（曖昧さ）と状態の微分化が混同されやすい。指標比較では BIC が最小かつ Silhouette が最大となるのが $K=2$ であり、本データでは $K=2$ が最も解釈可能である。ただし、特徴量を拡張した場合は $K > 2$ の潜在構造が現れる可能性があり、今後の課題とする。

3.2 二態の意味は「共同幸福」と「自己防衛」にする

なお、著者によるアンカーのスポットチェックでは、ラベル選定は作品解釈と整合し、修正は不要であった。

m/iso 空間では低 m・高 iso 近傍に一群、高 m・中 iso 近傍に一群が形成される。本研究では少量ラベル（アンカー）によりクラスタの意味を拘束（semantic constraint）し、解釈ラベルとして「自己防衛（Self_Protection）」と「共同幸福（Communal_Happiness）」を付与した。これはクラスタリング後の恣意的命名ではなく、アンカーによる意味固定を目的とした運用である。この拘束により、状態系列を本文（evidence）へ還元しやすい形で読める。

3.3 転調点は時系列上でスパイクとして現れる

boundaryness 時系列は複数の明確なスパイクを持つ。これは、どちらの状態にも属し切らない観測が局所的に生じていることを意味する。転調点は遷移回数と必ずしも一致せず、同じ章でも「決断としての明瞭な切替」は entropy を上げにくく、「葛藤・逡巡・曖昧な局面」は entropy を上げる。ゆえに転調点は心理的不確定性の極大として定義される。

3.3.1 ピーク3点の本文照合（例）

図3の boundaryness（事後確率エントロピー）時系列は、横軸が global_step（=イベント列の順序）であるため、ピーク位置を本文の引用（evidence）へ直接照合できる。本節では、entropy の局所最大のうち上位3点（global_step=20, 26, 39）を例として示す。

算出条件：K=2、 $\alpha=10$ 、reg=1e-3、diag=True（表3の設定）でアンカー付き EM を適用し、各点の事後確率 w_{tk} から $H_t = -\sum_k w_{tk} \log w_{tk}$ を計算。

表 3.3.1-1 entropy ピーク3点と本文引用（照合用）

Peak	global_step (x 位置)	章 (scene_id)	event	evidence (短引用)	boundaryne ss (p0/p1, H)
1	20	三、家 (3_home)	『ラッコの上 着』をからか われていると 語る	みんながぼく にあうとそれ を云う / ひや かすように云 うんだ	p0=0.751 p1=0.249 H=0.562
2	26	四、ケンタウ ル祭の夜 (4_festival_n ight)	怒って叫び返 し、理由を考 える	高く叫び返し ました / せわ しくいろいろ のことを考え	p0=0.137 p1=0.863 H=0.400
3	39	五、天気輪 (てんきり ん) の柱 (5_weather_ pillar)	列車と旅人を 想像して悲し みが強まる (推測)	汽車の音が聞 えてきました / かなしく なって	p0=0.751 p1=0.249 H=0.562

照合の見どころ（極簡易メモ）：

- global_step=20（家）：「みんなが…云う／ひやかす…」の社会的羞恥が明示され、前後（父の帰還を信じる→カムパネルラの配慮を回想）と合わせて“状態が揺れやすい”箇所としてピーク化しやすい。
- global_step=26（ケンタウル祭）：「叫び返しました／せわしく…考え」の反応・内省が、直前の攻撃（上着を投げつけられる）と直後の別の関心（星座早見）に挟まれ、自己防衛→再定位の境界として境界度が上がる。

- global_step=39（天気輪）：「汽車の音…／かなしくなって」で感情が変調し、直前の遠景（町の灯）から直後の“星空を温かな場所”と捉え直す流れの接続点としてピーク化しやすい。

【学術文献】

Montfort, N., & Pérez y Pérez, R. (2023). Computational Models for Understanding Narrative. *Revista de Comunicação e Linguagens*, (58), 97–117.

Meister, J. C. (2012). Crowd Sourcing Narrative Logic: Towards a Computational Narratology with CLEA. In *Proceedings of the Third Workshop on Computational Models of Narrative (CMN'12)*.

Bod, R., Fisseni, B., Kurji, A., & Lowe, B. (2012). Objectivity and Reproducibility of Proppian Narrative Annotations. In *Proceedings of CMN'12*.

Mani, I. (2013). Computational Modeling of Narrative. (NarrativeML 等を含む計算ナラトロジー総説).

Brown, M., Dobson, T., Grue, D., & Ruecker, S. (2013). Challenging New Views on Familiar Plotlines: Use of XML in a Scholarly Tool for Literary Pedagogy. *Literary & Linguistic Computing*, 28, 199–208.

Schwan, S., et al. (2019). Narrelations: Visualizing multiple levels of narrative. (HeureCLÉA/CATMA 系).

Reagan, A. J., Mitchell, L., Kiley, D., Danforth, C. M., & Dodds, P. S. (2016). The emotional arcs of stories are dominated by six basic shapes. *EPJ Data Science*.

Chen, A. T., Yoon, A., & Shaw, R. (2012). People, Places and Emotions: Visually Representing Historical Context in Oral Testimonies. In *CMN'12*.

Murai, H. (2018). Transitions of Plot Elements in a Japanese Detective Comic. (計量的プロット要素と遷移ネットワーク).

NeurIPS 2022. Revisiting Realistic Test-Time Training: Sequential Inference and Adaptation by Anchored Clustering.

【再現用リソース】

ginga-narrative-gmm（再現用リポジトリ）：GitHub リポジトリ（最終アクセス: 2026-01-02）

ginga_step_log_script.md（ステップログ可視化スクリプト）：GitHub ドキュメント（最終アクセス: 2026-01-02）

4. むすび

『銀河鉄道の夜』のイベント列は m/iso の 2 次元空間で二状態として安定に表現でき、事後確率エントロピーの極大点として転調点を定義することで、章構造と整合する形で状態遷移と転調の定量抽出が可能である。

5. 限界と今後の課題

本研究は、ループリック設計に内在する解釈バイアス、単一アノテータに起因する主観性、2 次元圧縮による情報落ちの影響を受ける。今後は多者アノテーションによる一致度評価、特徴量拡張（語り手距離、時間知覚、象徴密度等）、他作品への適用による外的妥当性検証が必要である。また、単一 AI による注釈は一貫性の利点がある一方、モデル固有のバイアスを含む可能性がある。今後は別モデルによる再注釈（同一プロンプト）を行い、アンカー一致度や転調点 Top-K の一致（Jaccard 等）でモデル間再現性を評価する。

付録 A : K-means と GMM/EM の違い（手法選択の補足）

本研究で GMM を採用する理由を、K-means との違いに基づき簡潔に整理する。

A.1 hard 割当と soft 割当

K-means は各点を必ず 1 クラスに確定割当する（hard）。一方 GMM は各点が各クラスに属する確率 w_{ik} （soft）を推定する。本研究の boundaryness は、この soft assignment の「曖昧さ」（エントロピー）を直接利用するため、GMM が自然である。

A.2 クラス形状の表現力

K-means は距離最小化に基づくため等方的（球状）な分割になりやすい。GMM は共分散 Σ_k を持ち、楕円形状（本実装では対角共分散も可）を表現できる。心理特徴量空間では状態の広がりが一様とは限らないため、GMM の柔軟性が有利となる。

A.3 最適化目的（SSE 最小化 vs 尤度最大化）

K-means は重心への二乗距離（SSE）を最小化する。一方 GMM は確率モデルの対数尤度を最大化するため、AIC/BIC など情報量基準で K を比較しやすい。なお、等方共分散かつ hard 割当へ近づく極限では、GMM の EM 更新は K-means（Lloyd 法）と密接に関係する。