

## # Pose-based Video Start Alignment on Mobile: Technical Note with Optional Audio-Assisted Refinement

### ## Pose系列による動画の開始点自動同期（端末内）：音声支援の粗密化・検証を含む技術ノート

> \*\*Open Technical Note (text-only).\*\* This document is shared to help practitioners reproduce and improve the method, and to quietly establish prior art for the described techniques.

\*\*Document license:\*\* CC BY 4.0 (attribution required).

\*\*Note:\*\* No patent license is granted by this document.

---

### ## Abstract (EN)

This technical note describes an on-device method to estimate the temporal offset between two practice videos containing the same or highly similar performance segment (e.g., dance/ballet), enabling synchronized side-by-side playback.

The method samples frames at a fixed FPS and extracts pose keypoints per frame using an on-device pose estimator. It covers both single-person and multi-person pose estimators; for multi-person outputs, a target person can be selected using detection score, average keypoint confidence, spatial consistency, track continuity, and/or region-of-interest constraints.

Keypoints are normalized into a pose feature vector by filtering low-confidence points, selecting a robust center (pelvis-first, then shoulders, then centroid fallback), and applying scale normalization (shoulder width, hip width, or maximum radius fallback). Temporal alignment is obtained by scanning candidate offsets within `±searchRangeSec` and computing an offset score as the average pose similarity across aligned frames. Per-frame similarity may be cosine similarity, and distance can be defined as  $1 - \cosSim$ . Frame pairs with fewer than a minimum number of commonly valid keypoints (shared-valid points) are excluded from scoring to improve robustness.

**\*\*Offset sign convention:\*\*** Let  $\Delta$  denote the shift applied to Video B relative to Video A when pairing frames for scoring. For a frame index  $t^*$  in A, the paired index in B is  $(t + \Delta)^*$ . Thus,  $\Delta > 0$  indicates that **B lags A** (B is evaluated later than A), and  $\Delta < 0$  indicates that **B leads A**.

The selected offset may be converted into a robust start point using the earliest contiguous above-threshold similarity segment (or equivalent window-based criteria), with staged relaxation retries if no such segment is found.

A confidence score determines acceptance:

```
```text
confidence = 0.5·bestSim + 0.3·margin + 0.2·validRate
```
```

```

Here, margin is the gap between best and runner-up offset scores, and validRate denotes a validity ratio computed over the matched segment/window (e.g., valid-frame ratio and/or valid-keypoint ratio). When distance-minimization is used internally, bestSim and runner-up values for margin can be computed on any monotonic transform consistent with the acceptance logic (e.g., converting distance to similarity), without loss of generality.

Optional hybrids are also described:

- Audio-assisted coarse candidate generation (top-K offsets with minimum separation) followed by pose refinement within local neighborhoods, and/or
- Post-pose audio verification or micro-refinement around the selected offset to improve precision.

Additionally, an optional pose-only phase acceleration is described (e.g., restricting search to early segments or high-motion windows, with fallback to full search).

**\*\*Keywords:\*\*** pose estimation, skeleton keypoints, temporal alignment, time offset estimation, sequence matching, cosine similarity, contiguous match, confidence scoring, coarse-to-fine, audio-to-pose, audio verification, micro-refinement, on-device, mobile, multi-pose

## ## 概要 (JP)

本技術ノートは、同一または高類似の演技区間を含む練習動画2本（例：バレエ／ダンス）を左右に並べて比較する際に、両動画の開始位置（時間オフセット）を端末内で自動推定し、同期再生に用いる方式をまとめたものです。実務者が再現・改善できる情報提供を目的に公開します。併せて、記載手法を先行技術（prior art）として記録する意図も含みます。

手法は、一定FPSでフレーム列を作成し、各フレームで骨格キーポイント（座標+信頼度）を推定する。單一人物推定に限らず、複数人物推定（Multi-Pose）の出力から対象人物を選別して用いる形態を含む（例：検出スコア、平均キーポイント信頼度、空間的一貫性、追跡の連續性、ROI制約など）。

低信頼キーポイントを除外した上で、中心（骨盤優先→肩→重心フォールバック）とスケール（肩幅／腰幅／最大半径フォールバック）により正規化した骨格ベクトル系列を構築する。次に、`±searchRangeSec` の範囲で時間オフセットを走査し、対応フレーム間の骨格ベクトル類似度（例：コサイン類似度）を平均集計したスコアで最良オフセットを選ぶ（距離は  $1 - \cos\text{Sim}$  と定義可能）。この際、両フレームで共通して有効なキーポイント数が所定数未満のフレーム対は、スコア計算対象から除外し頑健性を高める。

**\*\*符号規約：**  $\Delta$  を「Video A に対して Video B 側を評価する時間（フレーム）シフト」とし、Aのフレームtに対しBのフレーム(t+ $\Delta$ )を対応づける。 $\Delta > 0$  は「BがAより遅れている（BはAより後の時刻として評価される）」、 $\Delta < 0$  は「BがAより先行している」を意味する。

最良オフセットは、類似度が閾値以上となる連續区間（連續一致）を用いて頑健な開始点へ変換し、該当区間が見つからない場合は段階的に緩和して再試行する。

採否判定には以下の信頼度を用いる：

````text`

`confidence = 0.5×bestSim + 0.3×margin + 0.2×validRate`

...

\*本ドキュメントは CC BY 4.0 で提供します（クレジット表示で再利用可能）。本ドキュメントは特許ライセンスの許諾を意味しません。

---

## ## 1. Problem Setting / 課題設定

Two practice videos may contain the same or highly similar performance segment but start at different times. The goal is to estimate a temporal offset (start alignment) for synchronized side-by-side playback on mobile devices, robust to noise such as:

- different camera distance/zoom and minor viewpoint changes,
- partial occlusions and low-confidence keypoints,
- timing differences (pauses, warm-up motions, idle segments),
- pose-estimation noise and missing keypoints.

---

## ## 2. Inputs and Outputs / 入出力

### ### Inputs / 入力

- Video A, Video B
- Parameters (examples; not limiting):
  - `F`: sampling FPS (frames per second)
  - `searchRangeSec`: maximum offset search range ( $\pm$ , seconds)
  - `τ\_kp`: keypoint confidence threshold
  - `M\_shared`: minimum number of shared-valid keypoints for a frame-pair to be scored (e.g.,  $\geq 2$ ; implementation-defined)
  - `τ\_sim`: similarity threshold for contiguous match detection
  - `N`: minimum contiguous frames above threshold
  - `τ\_conf`: acceptance threshold for confidence
- Hybrid parameters (optional):
  - `k`: top-K audio candidates

- `minSepSec`: minimum separation between candidate peaks (seconds)
- `r\_frames` / `r\_sec`: local neighborhood radius around candidates (frames / seconds)
- `ε\_frames` / `ε\_sec`: micro-refinement radius for audio verification (frames / seconds)
- Optional runner-up separation (optional):
  - `δ\_excl\_frames` / `δ\_excl\_sec`: exclusion radius around  $\Delta^*$  when selecting runner-up

### ### Outputs / 出力

- `startA\_sec`, `startB\_sec`: aligned start timestamps for synchronized playback
- `confidence`: acceptance score (e.g., normalized to [0,1])
- Optional diagnostics:
  - `bestOffset\_frames`
  - `bestOffset\_sec`
  - `bestSim`
  - `secondBestSim`
  - `margin`
  - `validRate`

### ### Units / 単位 (重要)

Offsets and neighborhoods can be represented in either frames or seconds. In this note,  **$\Delta$  is expressed in frames by default** (integer frame shift), and its seconds representation is:

```
```text
Δ_sec = Δ / F
r_sec = r_frames / F
ε_sec = ε_frames / F
```
```

```

Equivalent embodiments may instead express all quantities directly in seconds.

---

## ## 3. Method / 手法

### ### 3.1 Frame Sampling / フレーム抽出

1. Decode each video into frames sampled at fixed FPS  $F$  and a fixed target size.
2. Frame index  $t$  corresponds to timestamp  $t/F$ .

### ### 3.2 Pose Estimation / 骨格推定 (Single-pose / Multi-pose)

For each frame, estimate pose keypoints with coordinates and confidence  $(x_i, y_i, s_i)$ .

- **Single-pose:** directly use the single-person output.
- **Multi-pose:** select a target person from multiple detected persons using one or more of:
  - detection score / instance score,
  - average keypoint confidence,
  - spatial consistency (e.g., torso center continuity, size constraints),
  - track continuity over time (simple tracking by nearest center / IoU / embedding),
  - region-of-interest constraints,
  - rejecting implausible skeleton geometry.

This note is not limited to any specific model or framework; any on-device pose estimator producing keypoints with confidence qualifies as an equivalent embodiment.

### ### 3.3 Pose Normalization & Vectorization / 正規化ベクトル化

Purpose: reduce variation due to camera translation and scale.

- Keypoint filtering:  $K_{valid}(t) = \{ i \mid s_i(t) \geq \tau_{kp} \}$
- Center selection (fallback): pelvis/hip center  $\rightarrow$  shoulder center  $\rightarrow$  centroid of valid keypoints
- Scale selection (fallback): shoulder width  $\rightarrow$  hip width  $\rightarrow$  max radius from center
- Normalized coordinates (for valid keypoints):
  - $p_i'(t) = ( (x_i(t) - c_x(t)) / \alpha(t), (y_i(t) - c_y(t)) / \alpha(t) )$
- Feature vector: concatenate normalized coordinates in a fixed keypoint order

Missing keypoints may be handled by masking (shared-valid only), zero-imputation with mask, or variable-length aggregation normalized by count.

### ### 3.4 Offset Scan & Score / オフセット走査とスコア

\*\*Offset sign convention:\*\* For frame index `t` in A, pair with `(t + Δ)` in B. Thus, \*\*Δ > 0\*\* means B lags A\*\*, and \*\*Δ < 0\*\* means B leads A\*\*.

Search range (frames):

- $\Delta \in [-\text{searchRangeSec}\cdot F, +\text{searchRangeSec}\cdot F]$

\*\*Discrete-range note (implementation-defined):\*\* since `searchRangeSec·F` may not be an integer, the conversion from seconds to frame bounds may apply rounding such as `round`, `floor`, or `ceil`, depending on implementation.

Shared-valid minimum:

- compute `sharedValidCount(t, Δ)` = number of keypoints valid in both frames
- exclude frame pairs with `sharedValidCount(t, Δ) < M\_shared`
- `M\_shared` is implementation-defined (e.g., ≥2) and tunable

Per-frame similarity:

- example: `sim(t, Δ) = cos( vA(t), vB(t+Δ) )`
- distance form: `dist(t, Δ) = 1 - sim(t, Δ)`

Offset score:

- `Score(Δ) = average\_t sim(t, Δ)` over included pairs (or minimize average distance)

Best and runner-up:

- $\Delta^* = \operatorname{argmax} \operatorname{Score}(\Delta)$
- $\Delta_2 = \operatorname{argmax}_{\{\Delta \neq \Delta^*\}} \operatorname{Score}(\Delta)$
- optional runner-up separation: exclude  $|\Delta - \Delta^*| \leq \delta_{\text{excl\_frames}}$  to avoid near-duplicate peaks

Optional (non-limiting): enforce a minimum overlap length and discard too-short overlaps.

### ### 3.5 Robust Start Point via Contiguous Match / 連続一致区間による開始点

- Find earliest `t0` such that  $\text{sim}(t, \Delta^*) \geq \tau_{\text{sim}}$  holds for at least `N` consecutive frames (valid frame pairs only).
- If not found, staged relaxation retries may lower  $\tau_{\text{sim}}$ , reduce `N`, switch to window-ratio criteria, switch to distance thresholds, etc.

Aligned start timestamps:

- Global search:
  - $\text{startA\_sec} = t0 / F$
  - $\text{startB\_sec} = (t0 + \Delta^*) / F$
- Windowed/local search (seeded neighborhoods, Phase1 restricted windows):
  - $\text{startA\_sec} = \text{baseA\_sec} + t0 / F$
  - $\text{startB\_sec} = \text{baseB\_sec} + (t0 + \Delta^*) / F$

### ### 3.6 Confidence & Acceptance / 信頼度と採否

Definitions:

- $\text{bestSim} = \text{Score}(\Delta^*)$
- $\text{secondBestSim} = \text{Score}(\Delta_2)$
- $\text{margin} = \text{bestSim} - \text{secondBestSim}$

Validity ratio `validRate` (over matched segment/window):

- $\text{validFrame}(t, \Delta) = 1$  if the pair is eligible and included (e.g.,  $\text{sharedValidCount} \geq M_{\text{shared}}$ ), else 0
- $\text{validFrameRate}$  = average of  $\text{validFrame}(t, \Delta^*)$
- optional: include validKeypointRate; validRate may denote valid-frame ratio and/or a combination

Confidence (example):

```
```text
confidence = 0.5 * bestSim + 0.3 * margin + 0.2 * validRate
````
```

Accept if  $\text{confidence} \geq \tau_{\text{conf}}$ , else reject.

Distance-minimization note:

- if computed via distance minimization internally, `bestSim` / `secondBestSim` for margin may be computed after any monotonic transform (e.g., distance → similarity)
  - values may be clipped/normalized to [0,1] prior to thresholding
- 

## ## 4. Optional Hybrids (Audio Assistance and/or Verification) / オプション：音声支援の粗候補生成および／または音声検証・微調整

### ### 4.1 Audio-based coarse candidate generation (optional) / 音声で粗候補生成（任意）

Compute audio similarity using any audio features (energy envelope, onset patterns, spectral features, fingerprint, cross-correlation, etc.). Generate top-K candidate offsets (in frames or seconds).

Optional robustness:

- enforce minimum separation between candidate peaks `minSepSec` to avoid near-duplicate candidates,
- keep both positive and negative candidates (respecting  $\Delta$  sign convention).

### ### 4.2 Pose refinement around candidates (optional) / Poseで候補近傍を精密化（任意）

For each candidate, refine only within a local neighborhood:

- $\Delta \in [\Delta_{c\_i} - r\_frames, \Delta_{c\_i} + r\_frames]$   
(or seconds:  $\Delta\_sec \in [\Delta_{c\_sec} - r\_sec, \Delta_{c\_sec} + r\_sec]$ )

Compute `Score( $\Delta$ )` using pose matching (Section 3), select best ` $\Delta^*$ ` across all neighborhoods, then apply contiguous-match and confidence.

### ### 4.3 Post-pose audio verification / micro refinement (optional) / Pose後の音声検証・微調整（任意）

After a pose-derived offset ` $\Delta^*$ ` is selected:

- evaluate audio similarity in a small neighborhood  $\Delta \in [\Delta^* - \varepsilon\_frames, \Delta^* + \varepsilon\_frames]$   
(or seconds)
- if a nearby offset yields better audio consistency (or passes an audio check), adjust

` $\Delta^*$ ` accordingly

- or reject if audio strongly contradicts the pose-derived alignment

---

## ## 5. Optional Pose-only Phase Search Acceleration (Phase1; example) / Poseのみの探索高速化（Phase1；例示）

To improve runtime on mobile devices, an optional pose-only acceleration phase can be used:

- **Early-segment restricted search:** search only within early segments (e.g., 0–10s, 10–20s) before expanding.
- **High-motion window selection:** detect windows with high motion magnitude and search only those windows. Examples of motion proxies include pose displacement norms across frames, simple frame-difference energy, optical-flow statistics (if available), and other lightweight motion proxies.
- **Fallback:** if no confident match is found in restricted windows, fall back to the full-range scan (Section 3.4).

This phase acceleration is orthogonal to audio-assisted hybrids and can be combined with them.

---

## ## 6. Failure Handling and Fallbacks / 失敗時のフォールバック

If pose alignment fails or `confidence <  $\tau_{\text{conf}}$ `:

- fall back to audio-only alignment (if available), and/or
- keep the pose offset but require audio verification before acceptance, and/or
- reject automatic alignment and allow manual adjustment.

---

## ## 7. Publication Metadata / 公開メタ情報

For a fixed release, record DOI and immutable URLs in external metadata (README / GitHub Release / Zenodo). If you also include them here, ensure you treat this file as

versioned content.

- Authors / 著者: Sparklyforce
- Version: v1.0.1
- First public timestamp (UTC, ISO 8601): YYYY-MM-DDThh:mm:ssZ
- Repository tag: v1.0.1
- Commit (Git SHA): <to be filled>
- GitHub Release (Permanent URL): <to be filled>
- Zenodo record (Immutable archive URL): <to be filled>
- Zenodo Version DOI: <to be filled>
- Zenodo Concept DOI (all versions): <to be filled>

---

## ## License / ライセンス

This document is licensed under \*\*Creative Commons Attribution 4.0 International (CC BY 4.0)\*\*.