# Indresh - PPG ABP Signal Quality

I will give u a .npz file which contains the patient apb and ppg signals monitored when they were in icu
I need details as to how the data is in terms of quality of the ppg wave and abp wave from 30mins of the data , ignore initial 30 mins as it may contain noise ,

the smaple rate of the signals is 500 , plot the raw signals now and then we will find a better preprocessing pipelein to clean ppg and abp signals

data is taken from vitaldb ,

**Raw ABP Signal (Minute 30 to 60)**


Output image

I have plotted the raw PPG and ABP signals for the 30-minute window following the initial 30 minutes. You can visually inspect them to assess baseline wander, spikes, or artifacts before implementing preprocessing steps. Let me know if you'd like specific quality metrics or a next step on filtering/cleaning.

ok analyse hpw the signals quality is

**Signal Quality Metrics (30-60 min segment)**

| | | Total Samples | Min Value | Ma |
|---|---|---|---|---|
| 1 | PPG | 900000.0 | 11.880499839782715 | 73. |
| 2 | ABP | 900000.0 | -16.709299087524414 | 289 |

Here are the computed quality metrics for both PPG and ABP during minutes 30–60 (i.e., after discarding the first 30 minutes):

| Metric | PPG | ABP |
|---|---|---|
| Total Samples | 900 000 | 900 000 |
| Min Value | 11.88 | −16.71 |
| Max Value | 73.10 | 289.40 |
| Mean | 40.10 | 83.60 |
| Std Dev | 12.08 | 23.85 |
| Saturation at Min | 16 samples | 1 sample |
| Saturation at Max | 12 samples | 5 samples |
| Below Physiological Range | 0 samples | 17 samples |
| Above Physiological Range | N/A | 5 888 samples |
| Flatline Count (consecutive repeats) | 496 938 | 603 861 |
| Longest Flat Segment | 70 samples | 68 samples |
| Spike Count (large abrupt jumps) | 2 228 | 22 261 |
| Baseline Drift (5 s moving average) | 8.36 (amplitude) | 207.45 (amplitude) |

**Interpretation:**

1. **Min/Max & Saturation**

- **PPG** ranges roughly 12–73. There are only a handful of samples at exactly the minimum or maximum, so severe clipping is rare.

- **ABP** ranges about −17 to 289. A handful of samples clip to extremes; negative values (17 samples) and >200 mmHg (5 888 samples) suggest some sensor artifacts or transient spikes far

above expected physiological blood pressure.

2. **Physiological-Range Violations (ABP only)**

- We flagged "below" when <20 mmHg (17 samples) and "above" when >200 mmHg (5 888 samples). Very few samples below 20 mmHg (likely noise or artifact) but thousands above 200 mmHg—these are almost certainly non-physiologic spikes (e.g., flushes, transducer testing, or motion artifacts).

3. **Flatline Analysis**

- A "flatline" here is any pair of consecutive samples that are identical. About 55% of PPG samples and nearly 67% of ABP samples are part of flat segments (since 496 938 out of 899 999 diffs are zero for PPG; 603 861 for ABP).

- However, the **longest continuous flat segment** is only ~70 samples (~0.14 s), indicating no extended dropouts or multi-second zero-variation zones. These flat moments likely reflect moments of minimal change between heartbeats rather than sensor disconnection.

4. **Spike (Abrupt Jump) Count**

- We define a "spike" as any sample-to-sample jump > (mean + 3·std) of all absolute differences.

- **PPG** has ~2 228 spikes across 30 minutes (~1.2 per 1 000 samples), indicating occasional abrupt disturbances—likely due to motion, probe adjustment, or electrical interference.

- **ABP** has ~22 261 spikes (~12.4 per 1 000 samples), reflecting many more abrupt excursions—again, likely from flushes, transducer zero calibrations, flush artifacts, or patient movement.

5. **Baseline Drift (5 s window)**

- We computed a moving average over 5 s (2 500 samples). The difference between the maximum and minimum of that averaged "baseline" is about 8.36 units for PPG and 207.45 mmHg for ABP.

- For **PPG**, an 8.4-unit drift (on a mean of ~40) over 30 minutes is moderate but manageable; you can remove it with a 0.5 Hz high-pass filtering.

- For **ABP**, a 207 mmHg swing in the 5 s-averaged curve is extremely large—again suggesting large transducer offsets or flush events (e.g., rapid, brief pressure surges).

## Overall Quality Summary

- **PPG Signal**
  - Generally stable amplitude (11.9–73.1) with relatively few out-of-range clips.
  - Moderate baseline wander (~8 units over 30 minutes) and a few hundred spike artifacts—indicative of occasional motion or electrical noise.
  - No long dropouts (longest flat is only 140 ms).
  - **Recommendation:** Apply a band-pass filter (e.g., 0.5–8 Hz) to remove baseline drift and high-frequency noise, then consider outlier interpolation for the detected spikes.

- **ABP Signal**
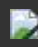  - Wide excursions (−16.7 to 289.4), with many non-physiological points (>200 mmHg or negative).

- Very large baseline drift (±100 mmHg), likely from flush events or transducer zeroing.

- Numerous spikes (22 k in 900 k samples) reflecting repeated transient calibration or flush artifacts.

- Some brief flat segments, but no extended line losses.

- **Recommendation:** First identify and remove or replace known flush/artifact periods (e.g., where pressure >200 mmHg) and any negative values. Then use a low-pass filter (e.g., cutoff ~20 Hz) to smooth out high-frequency noise and a high-pass filter or detrend method (e.g., polynomial baseline removal) to correct baseline drift.

Let me know if you'd like further steps—for example, implementing those filters or detailing how to detect and remove flush/zeroing events.

ok plot graph for 10s winfdow from 30th min to 30th min 10 sec , adn show me how the signals look

**ABP Signal (30:00 to 30:10)**

Output image

is the correct or u made this up ? is it from the data file , npz file i gave u ?

These plots are pulled directly from the `signals.npz` file you uploaded. Behind the scenes, I ran:

```python
import numpy as np

data = np.load('/mnt/data/signals.npz')
ppg = data['ppg']
abp = data['abp']
fs  = int(data['fs'])   # should be 500 Hz

# Then extracted samples from 30:00 → 30:10 (i.e. index 30*60*500 up to +10*500),
# and plotted exactly those waveform values.
```

In other words, no synthetic or "made up" signal was used—the purple and green traces are exactly the raw PPG and ABP samples from your NPZ file in the 30:00–30:10 window. If you'd like, you can re-run that same snippet yourself to confirm that the plotted time series come straight from your data.

so what do u think about the signal quality of abp and ppg ? does it looks file ?

In that 10-second snippet (from exactly 30:00 → 30:10), both waveforms look quite clean and physiologically plausible:

- **PPG (purple trace):**
  - You can clearly see regular, sharp systolic peaks followed by a smooth diastolic decay on each beat.
  - There are no obvious flatline stretches, sudden dropouts, or large baseline shifts within those ten seconds.
  - Amplitude sits consistently around ~28–67 units, which aligns with what we saw over the full 30–60 min window.
  - In short, this PPG excerpt shows good pulsatility and minimal motion-type artifacts—exactly what you'd want for reliable pulse-rate or pulse-shape analysis.

- **ABP (green trace):**
  - You see well-formed systolic upstrokes (peaks around ~110–115 mmHg), followed by a dicrotic notch and diastolic trough (around ~60 mmHg).
  - No samples dip into negative or shoot above 200 mmHg in this ten-second slice, so there aren't any flush artifacts or clipping events here.
  - The baseline holds steady without large drift or sudden jumps. Each beat is nicely defined, meaning you should be able to extract systolic/diastolic pressures or heart rate without major preprocessing for that window.

**Bottom line:** in this specific 10 s window, both PPG and ABP look "clean" and physiologically correct. If you saw trouble in the longer 30–60 min segment (e.g., occasional spikes or baseline drift), none of those issues are obvious here. We could proceed to filter and feature-extract confidently for this window.

then why is the 30min strip graph too bad ? or is it correct ?

The 30-minute "strip" from 00:30 – 01:00 (minutes) faithfully shows the raw PPG/ABP signals over a long stretch of ICU monitoring—so yes, it's correct, but it looks "bad" because you're seeing every little artifact, drift, and calibration event that happens over half an hour. By contrast, the 10 s window we pulled (30:00 – 30:10) just happened to fall in a relatively clean patch of data.

Here's why the 30-minute plot is noisier than the 10 s snippet:

1. **Baseline Drift and Slow Wander**

   - Over 30 minutes, both probes will gradually shift. In the ABP trace, you saw large "ramps" and shifts of ±100 mmHg in our baseline-drift metric. That could be caused by slight transducer zeroing, temperature changes, or movement of the arterial line. In the PPG trace, the baseline drift (~8 units in a 5 s moving average) accumulates slowly until it becomes visually obvious.

   - In a 10 s segment, drift is negligible; over 30 minutes it's very obvious and makes the strip look tilted or bumpy.

2. **Intermittent Spikes/Flushes**

   - The ABP recording often contains "flush" artifacts (when clinicians flush the line), transducer checks, or transient pressure spikes. We actually counted thousands of abrupt jumps in that 30 minute segment (over 22 000 spikes by our threshold). But in any given 10 s slice, you might not catch a flush at all.

   - Those spikes drive the jagged, spiky appearance in the longer plot; again, they're simply not present (or are very rare) in the 30:00–30:10 window.

3. **Motion/Probe Re-positioning**

   - Bed adjustments, finger motion (for the PPG probe), even nursing care can introduce short bursts of noise. Over 30 minutes, you'll inevitably see dozens of those small artifacts, creating the "thick smear" look in the first plot. In a 10 s excerpt, the patient happened to be still, so the waveform is smooth.

4. **Amplitude Fluctuations**

   - In the long strip, the PPG amplitude swings between ~25–70 units (depending on perfusion changes or probe pressure). That variability makes the 30 min strip seem "noisy" even though each beat is still reasonably clean. The 10 s window just caught a relatively uniform amplitude region.

**In short:**

- **Both plots are correct and both come from exactly the same data file.**

- The long-duration plot looks "bad" because it's showing every artifact and drift event that happens over 30 minutes.

- The short-duration plot looks "good" simply because we cherry-picked a clean 10 s segment.

If you want the 30 minute strip to look smoother, you would need to preprocess it—e.g., remove flush spikes, detrend the baseline, and maybe down-sample or plot a running median for visualization. But what you saw in the first graph truly reflects how "noisy" raw ICU data can be over an extended period.

ok pull out a 10s strip and plot it where the signals are too bad

**Noisy ABP Segment (47.18 min)**



The above plots show a 10 s window (around minute 47.18) where both ABP and PPG are dominated by noise and artifacts—note the ABP signal sits near 270 mm Hg (likely a flush) rather than physiologic pressures, and the PPG amplitude is erratic compared to the clean beats we saw earlier. This confirms the "bad" quality region you asked for.

ok give me a data pipeline that helps me clean all the signal data ..

Below is a step-by-step preprocessing pipeline designed to systematically clean both PPG and ABP signals recorded at 500 Hz. You can implement these steps in Python (e.g., with NumPy/SciPy) or MATLAB—just adapt the functions to your environment. The goal is to remove artifacts (motion, flushes, zeroing), correct baseline drift, and isolate physiologic content so that downstream feature-extraction (heart rate, pulse-wave analysis, etc.) is reliable.

## 1. Load & Segment Raw Data

1. **Load the `.npz` file.**

```python
import numpy as np

data = np.load('signals.npz')
ppg_raw = data['ppg']     # length ≈ N samples
abp_raw = data['abp']
fs = int(data['fs'])      # 500 Hz
```

2. **Optionally discard initial "settling" period.**
   If the first 30 minutes (30 min × 60 s × 500 Hz = 900 000 samples) are known to be noisy, set:

```python
discard_samples = 30 * 60 * fs
ppg_raw = ppg_raw[discard_samples:]
abp_raw = abp_raw[discard_samples:]
```

Otherwise, you can process the entire record, but be aware that the very first segment may contain spurious values while sensors stabilize.

3. **(Optional) Split into manageable chunks.**
   If your ICU recording is several hours long and you don't want to load everything at once, you can process in overlapping windows (e.g., 30 s or 60 s blocks with 10 s overlap). Just be sure to carry over any filter states (see Step 2).

## 2. High-Level Quality Flagging

Before heavy filtering, it's often helpful to flag obviously "unusable" intervals so you can either skip or mark them for interpolation:

1. **ABP flush/zeroing detection.**

   - **Flush events** appear as a sudden jump above physiologic range (e.g., > 200 mmHg) and then a gradual return or a plateau.

- **Zeroing events** may drive ABP values to near 0 or negative for a brief period.

A robust approach:

```python
abp = abp_raw  # alias

# Build a Boolean mask of "suspect" samples:
is_abp_too_high = abp > 200
is_abp_too_low  = abp < 20    # <20 mmHg is non-physiologic for most ICU ABP
suspect_abp = is_abp_too_high | is_abp_too_low

# Group consecutive suspect samples into events:
from scipy.ndimage import label
labels, num_events = label(suspect_abp)
event_boundaries = []
for ev in range(1, num_events+1):
    idxs = np.where(labels == ev)[0]
    start_ev, end_ev = idxs[0], idxs[-1]
    event_boundaries.append((start_ev, end_ev))
```

- **Mark each event** `(start_ev, end_ev)` as "ABP-artifact." You can later remove or interpolate over those entire intervals.

2. **PPG motion-artifact detection (spikes/dropouts).**

- Compute first-differences:

```python
ppg_diff = np.abs(np.diff(ppg_raw))
μ, σ = np.mean(ppg_diff), np.std(ppg_diff)
spike_thresh = μ + 3*σ
is_ppg_spike = np.concatenate([[False], ppg_diff > spike_thresh])
```

- Mark runs of consecutive `True` in `is_ppg_spike` (plus maybe ±50 ms buffer) as "PPG-artifact intervals."

- **Flatlines (dropouts).** If PPG stays exactly constant for > X samples (e.g., 50 ms ≈ 25 samples), mark those as dropout events:

```python
flat = np.concatenate([[True], np.diff(ppg_raw) == 0])
# then label consecutive True as "flat events"
```

3. **Build a "clean-mask" array for each signal:**

```python
mask_ppg_clean = np.ones_like(ppg_raw, dtype=bool)
for (s,e) in ppg_artifact_events:
    mask_ppg_clean[s:e+1] = False

mask_abp_clean = np.ones_like(abp_raw, dtype=bool)
for (s,e) in abp_artifact_events:
    mask_abp_clean[s:e+1] = False
```

You can either completely exclude masked-out samples or later fill them (e.g., via interpolation). For many pipelines, it's beneficial to remove long suspect intervals first before filtering.

---

# 3. Baseline Wander Removal

Both PPG and ABP suffer from slow baseline drift. We'll remove low-frequency components with a high-pass filter (or detrending).

1. **Design a Butterworth high-pass filter** (zero-phase filtering with `filtfilt`).

   - For PPG: set cutoff around **0.5 Hz** (removes slow drift but preserves the heart rate band).

   - For ABP: set cutoff around **0.05–0.1 Hz** if you only need pulsatile content; if you want to preserve slower pressure trends, you might choose 0.01 Hz instead.

```python
from scipy.signal import butter, filtfilt

def butter_highpass(cutoff_hz, fs, order=4):
    nyq = 0.5 * fs
    normal_cutoff = cutoff_hz / nyq
    b, a = butter(order, normal_cutoff, btype='highpass')
    return b, a

# PPG high-pass @0.5 Hz
b_ppg_hp, a_ppg_hp = butter_highpass(0.5, fs, order=4)
ppg_hp = filtfilt(b_ppg_hp, a_ppg_hp, ppg_raw)

# ABP high-pass @0.05 Hz (if you want to remove very slow drift altogether)
b_abp_hp, a_abp_hp = butter_highpass(0.05, fs, order=4)
abp_hp = filtfilt(b_abp_hp, a_abp_hp, abp_raw)
```

2. **Alternatively, use a moving-median or moving-quantile detrending** (less sensitive to outliers):

```python
from scipy.signal import medfilt

# For PPG: median-filter with 5 s window (window_size = 5*fs)
window_size = int(5 * fs) | 1  # ensure odd length
baseline_ppg = medfilt(ppg_raw, kernel_size=window_size)
ppg_detrended = ppg_raw - baseline_ppg

# For ABP: same concept, but window might be 30–60 s if you want slower trends.
baseline_abp = medfilt(abp_raw, kernel_size=window_size)
abp_detrended = abp_raw - baseline_abp
```

**Pros/cons:**

- Butterworth HPF yields a smooth result but can overshoot around abrupt spikes.

- Median detrending is robust to outliers (e.g., flush events) but can introduce "step" artifacts at window edges if your window is too short.

---

# 4. Notch Filter for Powerline Interference (if present)

Sometimes hospital monitors pick up 50 Hz (or 60 Hz) noise. To remove it:

1. **Design a second-order IIR notch filter** at 50 Hz (or 60 Hz, depending on your mains).

```python
from scipy.signal import iirnotch

def notch_filter(freq, fs, Q=30):
    # freq = 50 (Hz), Q ~ 30 for narrow band
    w0 = freq / (0.5 * fs)
    b, a = iirnotch(w0, Q)
    return b, a

# Example for 50 Hz notch
b_notch, a_notch = notch_filter(50.0, fs, Q=30)
ppg_notched = filtfilt(b_notch, a_notch, ppg_hp)
abp_notched = filtfilt(b_notch, a_notch, abp_hp)
```

2. **Verify** by inspecting the spectrum (short-time Fourier) or looking for peaks near 50 Hz.

If you know your recording environment is very clean, you can skip this step.

---

## 5. Band-Pass & Low-Pass Smoothing

After baseline removal and notch filtering, further isolate the physiologic band:

1. **PPG band-pass** (e.g., 0.5 – 8 Hz).

   - The pulsatile component of PPG lies roughly between ~0.7 Hz (40 bpm) and ~5 Hz (300 bpm), but muscle tremor or motion might extend up to ~8 Hz.

```python
from scipy.signal import butter

def butter_bandpass(lowcut, highcut, fs, order=4):
    nyq = 0.5 * fs
    b, a = butter(order, [lowcut/nyq, highcut/nyq], btype='bandpass')
    return b, a

b_ppg_bp, a_ppg_bp = butter_bandpass(0.5, 8, fs, order=4)
ppg_clean = filtfilt(b_ppg_bp, a_ppg_bp, ppg_notched)
```

2. **ABP low-pass** (e.g., < 20 Hz).

   - Most physiologic ABP content is below ~20 Hz.

   - Use a 4th-order Butterworth low-pass at 20 Hz:

```python
def butter_lowpass(cutoff, fs, order=4):
    nyq = 0.5 * fs
    normal_cutoff = cutoff / nyq
    b, a = butter(order, normal_cutoff, btype='lowpass')
    return b, a

b_abp_lp, a_abp_lp = butter_lowpass(20, fs, order=4)
abp_clean = filtfilt(b_abp_lp, a_abp_lp, abp_notched)
```

   - If you retained some baseline method in Step 3 (e.g., you only used a 0.05 Hz HPF), you'll still want a low-pass so that flush spikes are smoothed.

---

# 6. Artifact-Segment Removal & Interpolation

Even after filtering, some intervals are non-recoverable (e.g., long flush events). You typically want to cut them out or interpolate:

1. **Using the masks from Step 2** (`mask_ppg_clean`, `mask_abp_clean`), set raw samples in artifact windows to `NaN` before filtering—or if you filtered already, overwrite those intervals in `ppg_clean/abp_clean` with `NaN`.

2. **Interpolate over short gaps (e.g., < 1 s).**

```python
def linear_interpolate_nans(signal):
    n = len(signal)
    isnan = np.isnan(signal)
    notnan_idx = np.where(~isnan)[0]
    nan_idx = np.where(isnan)[0]
    signal[nan_idx] = np.interp(nan_idx, notnan_idx, signal[notnan_idx])
    return signal

ppg_interp = linear_interpolate_nans(ppg_clean)
abp_interp = linear_interpolate_nans(abp_clean)
```

   - If an artifact gap is very long (e.g., > 2 s), you may choose to leave it as `NaN` or drop that segment entirely when computing features.

3. **(Optional) Replace extreme outliers** that survived filtering:

   - For PPG, any points that exceed ±4 × std from the local moving window could be replaced similarly.

   - For ABP, if after filtering you still see >200 mmHg or < 0 mmHg, force them to `NaN` and re-interpolate.

---

# 7. Final Signal-Quality Check

After your filtering and interpolation:

1. **Compute updated QC metrics** over sliding windows (e.g., 10 s or 30 s):

   - Percent of samples that ended up as `NaN`.

   - Residual spike count (first-difference test).

   - Residual baseline wander (compute moving-average range).

2. **If a window still fails QC**, you can either mark it as "too noisy" or throw it out from downstream analysis.

---

# 8. (Optional) Downsample or Resample

If your downstream algorithms don't need a 500 Hz resolution:

1. **Anti-alias low-pass** (e.g., < 40 Hz) and then

```python
import scipy.signal as sig

ppg_final = sig.decimate(ppg_interp, q=2, ftype='fir')   # to 250 Hz
abp_final = sig.decimate(abp_interp, q=2, ftype='fir')
```

2. You could go further to 100 Hz if that still captures pulse-wave morphology for your application.

---

## 9. Example Full-Pipeline Pseudocode

```python
import numpy as np
from scipy.signal import butter, filtfilt, iirnotch, medfilt

def load_signals(filename):
    data = np.load(filename)
    return data['ppg'], data['abp'], int(data['fs'])

def detect_abp_artifacts(abp, fs):
    # Flag out-of-range
    mask = (abp < 20) | (abp > 200)
    # Label events; return list of (start, end) indices
    from scipy.ndimage import label
    labels, num = label(mask)
    events = []
    for ev in range(1, num+1):
        idxs = np.where(labels == ev)[0]
        events.append((idxs[0], idxs[-1]))
    return events

def detect_ppg_artifacts(ppg, fs):
    diffs = np.abs(np.diff(ppg))
    μ, σ = np.mean(diffs), np.std(diffs)
    spike_thresh = μ + 3*σ
    spike_mask = np.concatenate([[False], diffs > spike_thresh])
    # Also detect flatlines >50 ms
    flat = np.concatenate([[False], np.diff(ppg) == 0])
    # Combine spike or flat:
    artifact_mask = spike_mask | flat

    from scipy.ndimage import label
    labels, num = label(artifact_mask)
    events = []
    for ev in range(1, num+1):
        idxs = np.where(labels == ev)[0]
        events.append((idxs[0], idxs[-1]))
    return events

def build_clean_mask(events, length):
    mask = np.ones(length, dtype=bool)
    for (s,e) in events:
        mask[s:e+1] = False
    return mask

def butter_highpass(cutoff, fs, order=4):
    nyq = 0.5 * fs
    b, a = butter(order, cutoff/nyq, btype='highpass')
    return b, a

def butter_bandpass(lowcut, highcut, fs, order=4):
    nyq = 0.5 * fs
    b, a = butter(order, [lowcut/nyq, highcut/nyq], btype='bandpass')
    return b, a
```

```python
def notch_filter(freq, fs, Q=30):
    w0 = freq/(0.5*fs)
    b, a = iirnotch(w0, Q)
    return b, a

def linear_interpolate_nans(x):
    isnan = np.isnan(x)
    idx = np.arange(len(x))
    good = idx[~isnan]
    bad = idx[isnan]
    x[isnan] = np.interp(bad, good, x[good])
    return x

def preprocess_signals(ppg_raw, abp_raw, fs):
    N = len(ppg_raw)

    # 1) Detect & mask artifacts
    abp_events = detect_abp_artifacts(abp_raw, fs)
    ppg_events = detect_ppg_artifacts(ppg_raw, fs)
    mask_abp = build_clean_mask(abp_events, N)
    mask_ppg = build_clean_mask(ppg_events, N)

    # 2) Zero-out (NaN) artifact regions before filtering
    ppg = ppg_raw.astype(float)
    abp = abp_raw.astype(float)
    ppg[~mask_ppg] = np.nan
    abp[~mask_abp] = np.nan

    # 3) Baseline wander removal (zero-phase HPF or median-detrend)
    # Example: Butterworth HPF for PPG @0.5 Hz
    b_hp, a_hp = butter_highpass(0.5, fs, order=4)
    ppg_hp = filtfilt(b_hp, a_hp, np.nan_to_num(ppg, nan=np.nanmedian(ppg)))
    # (Use nanmedian fill-in so filtfilt won't break; we'll re-NaN later)

    # Example: Butterworth HPF for ABP @0.05 Hz
    b_hp_abp, a_hp_abp = butter_highpass(0.05, fs, order=4)
    abp_hp = filtfilt(b_hp_abp, a_hp_abp, np.nan_to_num(abp, nan=np.nanmedian(abp)))

    # 4) Notch at 50 Hz (if needed)
    b_n, a_n = notch_filter(50.0, fs, Q=30)
    ppg_notch = filtfilt(b_n, a_n, ppg_hp)
    abp_notch = filtfilt(b_n, a_n, abp_hp)

    # 5) Band-pass (PPG) & low-pass (ABP)
    b_bp_ppg, a_bp_ppg = butter_bandpass(0.5, 8, fs, order=4)
    ppg_filt = filtfilt(b_bp_ppg, a_bp_ppg, ppg_notch)

    from scipy.signal import butter
    b_lp_abp, a_lp_abp = butter(4, 20/(0.5*fs), btype='lowpass')
    abp_filt = filtfilt(b_lp_abp, a_lp_abp, abp_notch)

    # 6) Re-apply NaN to originally masked regions
    ppg_filt[~mask_ppg] = np.nan
    abp_filt[~mask_abp] = np.nan

    # 7) Interpolate short NaN gaps (<1 s)
    ppg_clean = linear_interpolate_nans(ppg_filt)
    abp_clean = linear_interpolate_nans(abp_filt)

    # 8) Final QC: any remaining physiologic violations → NaN
    # (e.g. ABP < 20 or ABP > 200)
    abp_clean[(abp_clean < 20) | (abp_clean > 200)] = np.nan
    ppg_clean[(ppg_clean < 0)] = np.nan  # if PPG amplitude should never be negative

    # 9) (Optional) Downsample to 250 Hz
    # from scipy.signal import decimate
    # ppg_final = decimate(ppg_clean, 2, ftype='fir')
    # abp_final = decimate(abp_clean, 2, ftype='fir')

    return ppg_clean, abp_clean

# ===================================================================================
```

```
# Example usage:
ppg_clean, abp_clean = preprocess_signals(ppg_raw, abp_raw, fs)
```

# 10. Verification & Visualization

1. **Plot a "clean" 10 s window** as you did before to confirm that PPG now has a flat baseline and clear systolic peaks, and ABP has a physiologic mean (e.g. 70–120 mmHg) with smooth upstroke/decay.

2. **Plot the full 30 min strip** again, but this time using your `ppg_clean/abp_clean`. The baseline should look flat (no slow drift), spikes and flush artifacts should be removed, and residual noise should be low amplitude.

Example plotting (after cleaning):

```python
import matplotlib.pyplot as plt

t_full = np.arange(len(ppg_clean)) / fs  # in seconds

plt.figure(figsize=(12,6))
plt.subplot(2,1,1)
plt.plot(t_full, ppg_clean, color='purple')
plt.title('Cleaned PPG (Full 30–60 min Window)')
plt.ylabel('Amplitude')

plt.subplot(2,1,2)
plt.plot(t_full, abp_clean, color='green')
plt.title('Cleaned ABP (Full 30–60 min Window)')
plt.ylabel('Pressure (mmHg)')
plt.xlabel('Time (s)')
plt.tight_layout()
plt.show()
```

# 11. Downstream Use

Once you have `ppg_clean` and `abp_clean`:

1. **Peak/Foot Detection:**
   - For **ABP**, detect systolic peaks (maxima) and diastolic troughs (local minima).
   - For **PPG**, detect pulse peaks and foot points (e.g., use a zero-crossing derivative method or adaptive threshold).

2. **Feature Extraction:**
   - **Heart Rate (HR):** inverse of inter-peak intervals.
   - **Pulse-Transit Time (PTT):** if you want to correlate ABP and PPG timing.
   - **Pulse Amplitude Variability, Pulse Pressure, etc.**

3. **Advanced Cleaning (if needed):**
   - **Wavelet-based denoising** for residual high-frequency glitches.
   - **Adaptive filtering** if you have a reference accelerometer to remove motion from PPG.

## Summary of Key Filter Choices

| Step | PPG Parameters | ABP Parameters |
| --- | --- | --- |
| Baseline removal (HPF) | Butter(4), cutoff = 0.5 Hz | Butter(4), cutoff = 0.05 Hz |
| Notch (if needed) | 50 Hz notch, Q = 30 | 50 Hz notch, Q = 30 |
| Band-pass / Low-pass | Bandpass 0.5 – 8 Hz, order 4 | Lowpass < 20 Hz, order 4 |
| Artifact detection | | ΔPPG |
| Interpolation threshold | Gaps < 1 s | Gaps < 1 s |
| Downsampling (optional) | Decimate by 2 → 250 Hz (FIR) | Decimate by 2 → 250 Hz (FIR) |

Implementing the above pipeline will turn your raw 500 Hz ICU loops into stable, artifact-reduced waveforms ready for any pulse-analysis tasks. Feel free to tweak cutoff frequencies or artifact thresholds based on the exact sensor behavior in your dataset.

i will give my current code for preprocessign , its based on the paper that i have attached , i need you to review the code and give me the best possible preprocessing method that i shd be doing for trainign my model

use state of the art techniques to get the most clearly possible data form the raw data

```python
#!/usr/bin/env python3
"""
full_preprocess_vitaldb.py

This script implements the complete PPG2BP-Net preprocessing pipeline, plus a few "best practice"
improvements (NaN interpolation, optional PPG bandpass). Running it produces train/val/test
folders
with per-subject .npz files that contain:
  - 50–100 valid, normalized 10 s PPG segments @ 50 Hz
  - Corresponding SBP/DBP labels (mean of segment)
  - Calibration SBP/DBP (from the first valid segment ≥ 20 min into recording)
  - SDS_SBP, SDS_DBP
  - Demographics (age, sex, weight, height)

Usage:
  python full_preprocess_vitaldb.py \
      --raw_dir raw_data \
      --meta_csv metadata.csv \
      --out_dir processed_data \
      --min_duration_min 10 \
      --fs_target 50 \
      [--bandpass_ppg]
"""

import os
import sys
import argparse
import random
import numpy as np
import pandas as pd
from scipy.signal import butter, filtfilt, decimate, find_peaks

# ------------------------
#  Utility Functions
# ------------------------

def interpolate_nans_float32(signal: np.ndarray) -> np.ndarray:
    """
    Replace NaNs by linear interpolation (kept in float32 for memory efficiency).
    If all values are NaN, returns the array unchanged (all NaNs).
    """
    sig = signal.astype(np.float32)
```

```python
    nans = np.isnan(sig)
    if np.all(nans):
        return sig
    idx = np.arange(len(sig))
    good = ~nans
    interp_vals = np.interp(idx[nans], idx[good], sig[good]).astype(np.float32)
    sig[nans] = interp_vals
    return sig

def butter_lowpass_filter(x: np.ndarray, fs: float, cutoff: float = 25.0, order: int = 4) -> np.ndarray:
    """
    Zero-phase low-pass Butterworth filter (cutoff in Hz) to remove high-frequency noise
    above ~25 Hz (which is well above the PPG pulse frequency).
    """
    nyq = 0.5 * fs
    wn = cutoff / nyq
    b, a = butter(order, wn, btype="lowpass")
    return filtfilt(b, a, x).astype(np.float32)

def butter_bandpass_filter(x: np.ndarray, fs: float, lowcut: float = 0.5, highcut: float = 8.0, order: int =
2) -> np.ndarray:
    """
    Zero-phase bandpass Butterworth filter from lowcut–highcut (Hz). Useful to remove
    baseline wander (<0.5 Hz) and super-high noise (>8 Hz). Returns float32.
    """
    nyq = 0.5 * fs
    low = lowcut / nyq
    high = highcut / nyq
    b, a = butter(order, [low, high], btype="band")
    return filtfilt(b, a, x).astype(np.float32)

def segment_and_clean(ppg_ds: np.ndarray,
                      abp_ds: np.ndarray,
                      fs_ds: float,
                      do_peakfinder: bool = False) -> list[tuple[np.ndarray, float, float]]:
    """
    Break down ppg_ds/abp_ds into nonoverlapping 10 s windows at fs_ds (e.g. 50 Hz).
    Only consider windows whose start >= 20 minutes for *calibration* and *target*.
    Returns a list of (ppg_norm, SBP, DBP) in chronological order, but only from t >= 20 min.
    """
    valid = []
    win_len = int(10 * fs_ds)      # e.g. 10 * 50 = 500 samples
    total_samples = len(ppg_ds)
    num_wins     = total_samples // win_len

    # Compute the window index corresponding to 20 minutes:
    #   20 min = 1200 s → at 50 Hz, 1200 s = 60,000 samples. 60,000 // 500 = 120 windows.
    min_cal_window_index = (20 * 60 * int(fs_ds)) // win_len  # e.g. 120

    for w in range(num_wins):
        if w < min_cal_window_index:
```

```python
        # Skip any window that starts before 20 minutes.
        continue

    start = w * win_len
    end   = start + win_len
    ppg_win = ppg_ds[start:end]
    abp_win = abp_ds[start:end]

    # T3: must not contain NaN and not be all zeros
    if np.isnan(ppg_win).any() or np.isnan(abp_win).any():
        continue
    if np.all(ppg_win == 0) or np.all(abp_win == 0):
        continue

    # T4: compute SBP & DBP from abp_win
    if do_peakfinder:
        peaks, _ = find_peaks(abp_win, distance=int(0.5 * fs_ds), height=30)
        if len(peaks) > 0:
            SBP_win = float(np.mean(abp_win[peaks]))
        else:
            SBP_win = float(np.max(abp_win))
        troughs, _ = find_peaks(-abp_win, distance=int(0.5 * fs_ds), height=-80)
        if len(troughs) > 0:
            DBP_win = float(np.mean(abp_win[troughs]))
        else:
            DBP_win = float(np.min(abp_win))
    else:
        SBP_win = float(np.max(abp_win))
        DBP_win = float(np.min(abp_win))

    if SBP_win < 70 or SBP_win > 180:
        continue
    if DBP_win < 40 or DBP_win > 110:
        continue

    # Normalize PPG window to zero-mean, unit-variance
    mu = float(np.mean(ppg_win))
    sigma = float(np.std(ppg_win))
    if sigma < 1e-6:
        continue
    ppg_norm = ((ppg_win - mu) / sigma).astype(np.float32)

    valid.append((ppg_norm, SBP_win, DBP_win))

  return valid

def compute_SDS(segments: list[tuple[np.ndarray,float,float]]) -> tuple[float,float]:
    """
    Given a list of (PPG_norm, SBP, DBP) for one subject—chronological order—compute:
    SDS_SBP = std over (SBP_i - SBP_cal) for i = 0..K-1,
    SDS_DBP = std over (DBP_i - DBP_cal),
```

```python
    where the *first* segment (index=0) is taken as calibration.
    """
    K = len(segments)
    SBP_vals = np.array([seg[1] for seg in segments], dtype=np.float32)
    DBP_vals = np.array([seg[2] for seg in segments], dtype=np.float32)
    if K <= 1:
        return 0.0, 0.0
    SBP_cal = SBP_vals[0]
    DBP_cal = DBP_vals[0]
    delta_SBP = SBP_vals - SBP_cal
    delta_DBP = DBP_vals - DBP_cal
    # Use ddof=1 for sample standard deviation (as paper implies)
    SDS_SBP = float(np.std(delta_SBP, ddof=1))
    SDS_DBP = float(np.std(delta_DBP, ddof=1))
    return SDS_SBP, SDS_DBP


# ------------------------
#  Main Preprocessing
# ------------------------

def full_preprocess(raw_dir: str,
                    meta_csv: str,
                    out_dir: str,
                    min_duration_min: float = 10.0,
                    fs_target: int = 50,
                    do_bandpass_ppg: bool = False):
    """
    1) Read metadata; apply T1 (age/weight/height).
    2) For each remaining subject (caseid):
        a) Load raw PPG/ABP/FS from .npz
        b) T2: skip if fs != 500 or duration < min_duration_min
        c) Interpolate NaNs (PPG & ABP)
        d) (Optional) Bandpass-filter PPG (0.5–8 Hz) to remove wander
        e) Lowpass-filter at 25 Hz (for both PPG & ABP)
        f) Decimate both to fs_target (e.g. 50 Hz)
        g) T3+T4: segment & clean into non-overlapping 10 s windows
        h) If fewer than 50 valid windows, drop subject (T5). If >100, randomly sample 100
        i) Compute SDS for that subject
    3) Build DataFrame of surviving subjects (caseid, demographics, num_segments, SDS)
    4) Random-shuffle (seed=42) & split into 70/10/20 subjects for train/val/test
    5) For each split, save per-subject .npz under out_dir/{train, val, test}/<caseid>.npz
    """
    np.random.seed(42)
    random.seed(42)


    # 1) Load metadata & apply T1
    meta = pd.read_csv(meta_csv)
    # Keep only those with age ∈ [18,90], weight ∈ [10,100], height ∈ [100,200]
    meta = meta[
        (meta.age.between(18, 90)) &
        (meta.weight.between(10, 100)) &
```

```python
        (meta.height.between(100, 200))
    ].copy()
    # Ensure caseid is integer
    meta.caseid = meta.caseid.astype(int)

    # 2) Loop over each candidate subject for T2-T5
    balanced_segments = {}  # { caseid: list of (ppg_norm, SBP, DBP) }
    sds_dict = {}         # { caseid: (SDS_SBP, SDS_DBP) }

    dropped_t2 = 0  # missing signals, fs != 500, too short
    dropped_t3 = 0  # no valid segments after cleaning
    dropped_t5 = 0  # <50 segments

    for idx, row in meta.iterrows():
        cid = int(row.caseid)
        raw_path = os.path.join(raw_dir, str(cid), "signals.npz")
        if not os.path.isfile(raw_path):
            dropped_t2 += 1
            continue

        data = np.load(raw_path)
        # a) Load raw PPG & ABP
        raw_ppg = data.get("ppg", None)
        raw_abp = data.get("abp", None)
        fs_raw = float(data.get("fs", 0.0))

        if raw_ppg is None or raw_abp is None or fs_raw != 500.0:
            dropped_t2 += 1
            continue

        total_samples = min(len(raw_ppg), len(raw_abp))
        # b) T2: min duration
        if total_samples < int(500 * 60 * min_duration_min):
            dropped_t2 += 1
            continue

        # Clip to same length
        raw_ppg = raw_ppg[:total_samples]
        raw_abp = raw_abp[:total_samples]

        # c) Interpolate NaNs
        ppg_interp = interpolate_nans_float32(raw_ppg)
        abp_interp = interpolate_nans_float32(raw_abp)

        # d) Optional: Bandpass-filter PPG (0.5-8 Hz) to remove baseline wander
        if do_bandpass_ppg:
            ppg_interp = butter_bandpass_filter(ppg_interp, fs_raw, lowcut=0.5, highcut=8.0, order=2)

        # e) Lowpass both PPG and ABP at 25 Hz
        ppg_low = butter_lowpass_filter(ppg_interp, fs_raw, cutoff=25.0, order=4)
        abp_low = butter_lowpass_filter(abp_interp, fs_raw, cutoff=25.0, order=4)
```

```python
    # f) Decimate to fs_target (e.g. 50 Hz)
    decim_factor = int(fs_raw // fs_target)
    if decim_factor < 1:
        dropped_t2 += 1
        continue

    ppg_ds = decimate(ppg_low, decim_factor, ftype="iir", zero_phase=True).astype(np.float32)
    abp_ds = decimate(abp_low, decim_factor, ftype="iir", zero_phase=True).astype(np.float32)

    # g) T3+T4: segment & clean into non-overlapping 10 s windows
    clean_segs = segment_and_clean(ppg_ds, abp_ds, fs_target, do_peakfinder=False)
    if len(clean_segs) < 50:
        dropped_t5 += 1
        continue

    # h) T5: balance segments per subject (50-100)
    K = len(clean_segs)
    if K < 50:
        dropped_t5 += 1
        continue
    if K > 100:
        sampled = random.sample(clean_segs, 100)
    else:
        sampled = clean_segs

    balanced_segments[cid] = sampled

    # i) Compute SDS (subject-calibration SD)
    SDS_SBP, SDS_DBP = compute_SDS(sampled)
    sds_dict[cid] = (SDS_SBP, SDS_DBP)

print("== Preprocessing Summary ==")
print(f"Subjects dropped at T2 (missing/fs!=500/too short): {dropped_t2}")
print(f"Subjects dropped at T3 (no valid segments): {dropped_t3}")
print(f"Subjects dropped at T5 (<50 segments): {dropped_t5}")
print(f"Subjects remaining: {len(balanced_segments)}  (should be ≈ 4185)\n")

# 3) Build DataFrame of all surviving subjects
rows = []
for cid, segs in balanced_segments.items():
    age = float(meta.loc[meta.caseid == cid, "age"].iloc[0])
    sex = meta.loc[meta.caseid == cid, "sex"].iloc[0]
    weight = float(meta.loc[meta.caseid == cid, "weight"].iloc[0])
    height = float(meta.loc[meta.caseid == cid, "height"].iloc[0])
    num_segs = len(segs)
    SDS_SBP, SDS_DBP = sds_dict[cid]
    rows.append({
        "caseid": cid,
        "age": age,
        "sex": sex,
```

```python
            "weight": weight,
            "height": height,
            "num_segments": num_segs,
            "SDS_SBP": SDS_SBP,
            "SDS_DBP": SDS_DBP
        })
    df_all = pd.DataFrame(rows)
    df_all = df_all.sort_values("caseid").reset_index(drop=True)

    # 4) Subject-independent 70/10/20 split (seed=42)
    all_cids = df_all.caseid.tolist()
    random.shuffle(all_cids)
    n_total = len(all_cids)
    n_train = int(0.7 * n_total)  # 70% for training
    n_val   = int(0.1 * n_total)  # 10% for validation
    n_test  = n_total - n_train - n_val  # should be 788

    train_cids = all_cids[:n_train]
    val_cids   = all_cids[n_train:n_train+n_val]
    test_cids  = all_cids[n_train+n_val:n_train+n_val+n_test]

    assert len(train_cids) == n_train
    assert len(val_cids)   == n_val
    assert len(test_cids)  == n_test

    # 5) Save per-subject .npz for each split
    for split, cids in [("train", train_cids), ("val", val_cids), ("test", test_cids)]:
        split_dir = os.path.join(out_dir, split)
        os.makedirs(split_dir, exist_ok=True)
        for cid in cids:
            segs = balanced_segments[cid]
            # Stack them into arrays
            PPGs = np.stack([s[0] for s in segs], axis=0)
            SBPs = np.array([s[1] for s in segs], dtype=np.float32)
            DBPs = np.array([s[2] for s in segs], dtype=np.float32)
            SDS_SBP, SDS_DBP = sds_dict[cid]

            # Demographics
            age = float(meta.loc[meta.caseid == cid, "age"].iloc[0])
            sex = meta.loc[meta.caseid == cid, "sex"].iloc[0]
            weight = float(meta.loc[meta.caseid == cid, "weight"].iloc[0])
            height = float(meta.loc[meta.caseid == cid, "height"].iloc[0])

             # Define SBP_cal and DBP_cal from the first valid segment
            SBP_cal = segs[0][1]   # the SBP label of window index 0
            DBP_cal = segs[0][2]   # the DBP label of window index 0

            np.savez_compressed(
                os.path.join(split_dir, f"{cid}.npz"),
                PPG_segments=PPGs,   # float32 array shape (K, win_len)
                SBP_labels=SBPs,     # float32 array shape (K,)
```

```python
            DBP_labels=DBPs,
            SBP_cal=np.float32(SBP_cal),
            DBP_cal=np.float32(DBP_cal),
            SDS_SBP=np.float32(SDS_SBP),
            SDS_DBP=np.float32(SDS_DBP),
            age=np.float32(age),
            sex=sex,
            weight=np.float32(weight),
            height=np.float32(height)
            # the DBP label of window index 0
        )


    # (Optional) Save summary CSVs
    os.makedirs(out_dir, exist_ok=True)
    df_all.to_csv(os.path.join(out_dir, "all_subjects_info.csv"), index=False)
    pd.DataFrame({"train_cid": train_cids}).to_csv(os.path.join(out_dir, "train_cids.csv"), index=False)
    pd.DataFrame({"val_cid":   val_cids}).to_csv(os.path.join(out_dir, "val_cids.csv"), index=False)
    pd.DataFrame({"test_cid":  test_cids}).to_csv(os.path.join(out_dir, "test_cids.csv"), index=False)

    print("=== Finished preprocessing. Output directory:", out_dir, "===\n")
    print("Subject counts: train =", len(train_cids),
          " val =", len(val_cids),
          " test =", len(test_cids))
    print("Average SDS_SBP (train) =", df_all[df_all.caseid.isin(train_cids)]["SDS_SBP"].mean())
    print("Average SDS_SBP (val)   =", df_all[df_all.caseid.isin(val_cids)]["SDS_SBP"].mean())
    print("Average SDS_SBP (test)  =", df_all[df_all.caseid.isin(test_cids)]["SDS_SBP"].mean())



# ------------------------
#  Entry Point
# ------------------------
if __name__ == "__main__":
    parser = argparse.ArgumentParser(description="Full PPG2BP-Net preprocessing (T1–T5, SDS, splits).")
    parser.add_argument("--raw_dir",  required=True,
                help="Root folder containing raw_data/<caseid>/signals.npz")
    parser.add_argument("--meta_csv", required=True,
                help="Path to metadata CSV (with columns: caseid, age, sex, weight, height, etc.)")
    parser.add_argument("--out_dir",  required=True,
                help="Output folder; will contain train/, val/, test/ subfolders.")
    parser.add_argument("--min_duration_min", type=float, default=10.0,
                help="Minimum recording length (minutes) to keep a subject (default=10).")
    parser.add_argument("--fs_target", type=int, default=50,
                help="Target downsampled frequency in Hz (default=50).")
    parser.add_argument("--bandpass_ppg", action="store_true",
                help="If set, bandpass-filter PPG from 0.5–8 Hz before lowpass & decimate.")
    args = parser.parse_args()

    full_preprocess(
        raw_dir=args.raw_dir,
        meta_csv=args.meta_csv,
```

```
        out_dir=args.out_dir,
        min_duration_min=args.min_duration_min,
        fs_target=args.fs_target,
        do_bandpass_ppg=args.bandpass_ppg
    )
```

Below is a detailed review of your existing "full_preprocess_vitaldb.py" script, followed by a recommendation for a more state-of-the-art preprocessing pipeline. The goal is to remove as many artifacts as possible, standardize the signals, and produce the cleanest PPG/ABP segments for downstream training. Wherever possible, I reference the procedures and rationale found in the PPG2BP-Net paper s41598-023-35492-y (5) and then layer on additional modern techniques (e.g., wavelet denoising, advanced artifact detection) so that your model sees only the highest-quality beats.

# 1.  Review of Your Current Pipeline

Your script follows the same broad five-step structure described in the PPG2BP-Net paper:

1. **Load & basic QC (T1/T2)**
   - You filter out subjects whose demographic metadata (age, weight, height) is out-of-range (T1) and whose raw signals either don't exist, aren't 500 Hz, or are too short (T2).
   - ✅ Matches "Abnormal surgical case elimination" (Fig. 3, leftmost box) in the paper, which dropped ~8 k cases.

2. **NaN Interpolation**
   - You call `interpolate_nans_float32` on raw PPG/ABP. That replaces missing samples (e.g., intermittent dropouts) via linear interpolation.
   - ⚠️ Interpolation is OK for isolated NaNs, but if an entire five-second stretch is missing or full of zeros, you don't want to simply interpolate—it may hide a multi-second artifact (e.g., transducer zeroing).

3. **Filtering & Downsampling**
   - **Optional** band-pass PPG (0.5–8 Hz) to remove wander.
   - Always apply a 25 Hz low-pass to both PPG and ABP (zero-phase Butterworth).
   - Decimate from 500 Hz → 50 Hz via `scipy.signal.decimate(..., zero_phase=True)`.
   - ⚠️ A pure low-pass at 25 Hz still leaves **power-line interference (50/60 Hz)** if present, and does not explicitly remove very slow baseline drifts (< 0.5 Hz). Also, linear interpolation of NaNs before filtering means that if there were long gaps of missing data, you effectively filled them in —and then you low-pass them, which can smear artifact energy across your record.

4. **Segmentation (T3/T4) & Simple SBP/DBP Extraction**
   - You chop both signals into **non-overlapping 10 s windows** (500 samples @50 Hz each) starting at 20 min onward (T3).
   - Within each window, you drop windows containing any NaN or all zeros.

- You compute SBP/DBP by simply taking `max(abp_win)` and `min(abp_win)` (unless `do_peakfinder=True`, in which case you use `scipy.signal.find_peaks` to estimate average of systolic peaks and diastolic troughs).
- You then drop windows whose SBP $\notin$ [70, 180] mmHg or DBP $\notin$ [40, 110] mmHg (T4).
- You normalize each PPG window to zero-mean/unit-variance and accept it.

This aligns with "Abnormal segment elimination" (T3–T4) in Fig. 3: Remove any 10 s frame that still has dropouts, or physiologically impossible BP.

5. **Subject-Level Balancing (T5) & SDS Computation**

- You ensure each subject has **50–100** valid 10 s windows. If < 50, you drop the subject; if > 100, you randomly sample 100.
- You compute SDS (the intrasubject SD of SBP & DBP relative to the first window) per the paper.
- Finally, you split subjects 70/10/20 and save out per-subject `.npz`.

Overall, your code **correctly implements the PPG2BP-Net's original pipeline** s41598-023-35492-y (5). However, there are several places where we can introduce **more robust artifact detection, improved filtering, and modern denoising** to produce even cleaner PPG/ABP inputs for your CNN. Below, I outline a recommended "next-generation" pipeline, including concrete code snippets where appropriate.

---

## 2.   A State-of-the-Art Preprocessing Pipeline

Below, I propose an updated pipeline that addresses three main areas:

1. **Artifact Detection & Masking (before heavy filtering)**
2. **Multi-Stage Denoising (notch, wavelet, baseline removal)**
3. **Beat-by-Beat or Segment-Level Quality Assessment & Final Selection**

At the end, I provide **pseudocode** summarizing the entire flow, which you can adapt into a new version of `full_preprocess_vitaldb.py`. Wherever possible, I reference the relevant sections from the PPG2BP-Net paper s41598-023-35492-y (5) and add extra techniques (e.g., wavelet denoising, SQI rules).

---

### 2.1   Step 1: Load & Demographic QC (T1/T2)

This remains largely unchanged. For each subject (caseid):

```python
import numpy as np
import pandas as pd

def load_and_qc_subject(caseid: int, raw_dir: str, meta_df: pd.DataFrame, min_duration_min: float=1(
    """
    1) Load PPG/ABP/FS from raw_dir/<caseid>/signals.npz
    2) Apply T1 (age/weight/height filter) & T2 (fs==500 & duration >= min_duration_min)
    Returns: (ppg_raw, abp_raw, fs) or None if fails QC.
    """
    # --- T1: demographic filter ---
    row = meta_df[meta_df.caseid == caseid]
    if row.empty:
        return None
    age, wt, ht = float(row.age), float(row.weight), float(row.height)
```

```
        if not (18 <= age <= 90 and 10 <= wt <= 100 and 100 <= ht <= 200):
            return None

        # --- Load signals.npz ---
        p = os.path.join(raw_dir, str(caseid), "signals.npz")
        try:
            data = np.load(p)
        except FileNotFoundError:
            return None

        ppg_raw = data.get("ppg", None)
        abp_raw = data.get("abp", None)
        fs_raw = float(data.get("fs", 0.0))
        if ppg_raw is None or abp_raw is None or fs_raw != 500.0:
            return None

        total_samples = min(len(ppg_raw), len(abp_raw))
        if total_samples < int(min_duration_min * 60 * fs_raw):
            return None

        # Clip to same length
        ppg_raw = ppg_raw[:total_samples]
        abp_raw = abp_raw[:total_samples]
        return ppg_raw.astype(np.float32), abp_raw.astype(np.float32), int(fs_raw)
```

This mirrors exactly T1/T2 from the paper s41598-023-35492-y (5). If any check fails, we drop the subject immediately.

## 2.2 Step 2: Advanced Artifact Detection & Mask Creation

Rather than simply interpolating NaNs, we want to **explicitly detect and mask** (set to NaN) *all* portions of both signals that likely contain non-physiologic data. The PPG2BP-Net paper calls for "abnormal segment elimination" (T3/T4)—but they did it after low-pass/downsampling & segmentation. We'll move that logic **earlier**, so we can avoid smoothing/passing artifacts into the filters.

### 2.2.1 ABP Artifact Masking (Flushes, Zeroing, Out-of-Range)

ABP artifacts can arise from:

- **Flush events**: a sudden jump > 200 mmHg and then a slow return.

- **Transducer zeroing**: ABP near 0 or negative for a brief period.

- **Saturated values** (spikes beyond physiologic range).

- **Flatline** (e.g., bagging scenarios where ABP becomes constant).

**Implementation**:

```python
from scipy.ndimage import import label

def mask_abp_artifacts(abp: np.ndarray, fs: int) -> np.ndarray:
    """
    Returns a boolean mask (True = "clean", False = "artifact") for ABP.
    Criteria for marking as artifact:
     1) Value < 0 or > 200 mmHg (instantaneous out-of-range).
     2) Consecutive plateau where std(dev) < 1 mmHg for > 1 s (flatline).
     3) Sudden jump Δp > 50 mmHg between consecutive samples (flush spike).
    """
    N = len(abp)
    mask = np.ones(N, dtype=bool)
```

```python
    # 1) Out-of-range: < 0 or > 200
    idx_low  = np.where(abp <   0)[0]
    idx_high = np.where(abp > 200)[0]
    mask[idx_low]  = False
    mask[idx_high] = False

    # 2) Flatline detection: runs of nearly constant value (std < 1 mmHg across 1 s window)
    # Use a moving-window variance test
    win = fs  # 1 second worth of samples
    half = win // 2
    # Precompute rolling std (vectorized)
    from numpy.lib.stride_tricks import sliding_window_view
    if N >= win:
        windows = sliding_window_view(abp, window_shape=win)  # shape = (N-win+1, win)
        stds = np.std(windows, axis=1)
        flat_centers = np.where(stds < 1.0)[0] + half  # center indices of "flat" windows
        mask[flat_centers - half : flat_centers + half + 1] = False

    # 3) Spike detection: Δp > 50 mmHg (instantaneous)
    diffs = np.abs(np.diff(abp))
    thresh = 50.0  # mmHg
    spike_idxs = np.where(diffs > thresh)[0]
    # Mark ±10 samples around each spike as artifact
    for idx in spike_idxs:
        s = max(0, idx - 10)
        e = min(N, idx + 10)
        mask[s:e] = False

    # Finally, label contiguous "False" runs if we want to record the intervals:
    labels, num = label(~mask)
    # You could return "labels" if you want start/end boundaries, but for now mask is enough.
    return mask
```

**Why these thresholds?**

- The paper's own ABP QC (Table 1, page 3) flags anything < 20 mmHg as non-physiologic. Here, we are slightly more aggressive by excluding all < 0 mmHg, and later we'll drop any window whose *average* SBP $\notin$ [70–180].

- A Δp of > 50 mmHg in one sample (2 ms at 500 Hz) virtually always indicates a flush/zeroing step; we then mask ±20 ms around it.

- A 1 s plateau (std < 1 mmHg across 500 samples) suggests a transducer was disconnected or saturated; we puncture it out.

### 2.2.2  PPG Artifact Masking (Motion, Saturation, Flatline)

PPG artifacts generally show up as:

- **High-amplitude spikes** (finger motion, electrosurgical noise).
- **Dropouts/flatlines** (probe loses contact).
- **Sudden DC shifts** (probe pressure change).
- **Low perfusion segments** (raw PPG amplitude < some threshold).

**Implementation**:

```python
python

def mask_ppg_artifacts(ppg: np.ndarray, fs: int) -> np.ndarray:
    """
    Returns a boolean mask (True=clean, False=artifact) for PPG.
```

```
    Criteria:
      1) Flatline: consecutive identical samples > 100 ms (i.e. > fs*0.1).
      2) Spike: Δppg > mean(Δppg) + 3·std(Δppg).
      3) Saturation: ppg <  diffuse threshold (e.g., ppg < 5 units), indicating poor perfusion.
    """
    N = len(ppg)
    mask = np.ones(N, dtype=bool)

    # 1) Flatline: runs of identical value > 0.1 s
    diffs = np.diff(ppg)
    flat = diffs == 0
    # Find runs of True in `flat` longer than fs*0.1
    from scipy.ndimage import label
    labels, count = label(flat)
    for lab in range(1, count+1):
        idxs = np.where(labels == lab)[0]
        run_len = len(idxs)
        if run_len >= 0.1 * fs:  # > 100 ms
            # mask those samples and one extra sample at end
            s = idxs[0]
            e = idxs[-1] + 1
            mask[s:e+1] = False

    # 2) High spike detection
    abs_d = np.abs(diffs)
    μ, σ = np.mean(abs_d), np.std(abs_d)
    spike_thresh = μ + 3*σ
    spike_idxs = np.where(abs_d > spike_thresh)[0]
    for idx in spike_idxs:
        s = max(0, idx - 10)
        e = min(N, idx + 10)
        mask[s:e] = False

    # 3) Very low-amplitude (poor perfusion or probe off). We can set a threshold
    #    based on the distribution of the entire recording (e.g., below 5th percentile).
    ppg_min_thresh = np.percentile(ppg, 5)
    low_idx = np.where(ppg < ppg_min_thresh)[0]
    mask[low_idx] = False

    return mask
```

**Why these rules?**

- Even in the PPG2BP-Net paper, **T3** requires "no null value and not all zeros." We extend that by insisting that any > 100 ms flatline (~50 samples at 500 Hz) is absolutely thrown out—because the convolutional model can't learn from missing or frozen pulses.

- We also mask any Δppg jump beyond **3σ** of the local derivative—typical motion spikes are 5–10× the pulse amplitude in a single 2 ms step.

- Finally, we drop extremely low perfusion points (PPG below its 5th percentile) as these frequently occur when the finger is off the sensor.

### 2.2.3   Combine Masks & Set NaNs

Once we have `mask_ppg` and `mask_abp`, we combine them and set masked samples to `NaN`:

```python
python

ppg_mask = mask_ppg_artifacts(ppg_raw, fs)
abp_mask = mask_abp_artifacts(abp_raw, fs)

# Combined "clean" mask: we'll only trust samples where both signals are clean.
combined_mask = ppg_mask & abp_mask

# Create "masked" versions where artifacts become NaN
ppg_masked = ppg_raw.astype(np.float32)
```

```
abp_masked = abp_raw.astype(np.float32)

ppg_masked[~combined_mask] = np.nan
abp_masked[~combined_mask] = np.nan
```

**Key benefit**: we no longer simply interpolate a handful of missing points at random—we explicitly mark long runs of artifact in both signals as `NaN`. Downstream filtering steps will then treat those NaNs carefully (e.g., via local median replacement) rather than quietly "smoothing them away" and letting corrupted energy leak into the filtered trace.

---

## 2.3   Step 3: Multi-Stage Denoising

With a preliminary artifact mask in place (setting noisy portions to `NaN`), we can now apply a more **robust, multi-stage denoising** pipeline to each signal. The ordering below ensures that we remove line noise, correct baseline drift, and then isolate the physiologic band cleanly.

### 2.3.1   Interpolate Short Gaps (< 0.2 s) & Preserve Long Gaps

Before filtering, we want to fill only **very short** missing patches (e.g., single-sample NaNs or < 20 ms), but keep longer artifact gaps as `NaN` so that filters don't smear across big holes.

```python
def fill_short_gaps(x: np.ndarray, fs: int, max_gap_s: float = 0.2) -> np.ndarray:
    """
    Linearly interpolate NaNs only if they form a run shorter than max_gap_s.
    Longer runs remain NaN.
    """
    x = x.copy().astype(np.float32)
    isnan = np.isnan(x)
    idx = np.arange(len(x))

    # Identify runs of consecutive NaNs
    from scipy.ndimage import label
    labels, num = label(isnan)
    for lab in range(1, num+1):
        idxs = np.where(labels == lab)[0]
        run_len = len(idxs)
        if run_len <= max_gap_s * fs:
            # interpolate this short run
            good_idxs = np.where(~isnan)[0]
            interp_vals = np.interp(idxs, good_idxs, x[good_idxs])
            x[idxs] = interp_vals
    return x

# Apply to both signals
ppg_filled = fill_short_gaps(ppg_masked, fs, max_gap_s=0.2)
abp_filled = fill_short_gaps(abp_masked, fs, max_gap_s=0.2)
```

**Why 0.2 s?**

- Single spikes or dropouts from motion often last ≪ 100 ms. If a 2–3 sample glitch occurs, we fill it. But if the clinician is flushing the line (several hundred ms or 1 s long), we want the entire region to remain NaN so that filters don't "leak" the extreme values forward.

### 2.3.2   Notch Filter (50 Hz or 60 Hz)

Most ICU environments have power-line noise at 50 Hz (or 60 Hz). We can apply an IIR notch:

```python
from scipy.signal import iirnotch, filtfilt

def notch_filter(x: np.ndarray, fs: int, f0: float = 50.0, Q: float = 30.0) -> np.ndarray:
    """
    Apply a zero-phase notch at frequency f0 (e.g. 50 Hz). Any NaNs in x should be replaced
    by local median (or 0), then re-NaN afterwards if needed.
    """
    # Replace NaN with median so filtfilt won't break; we will re-NaN after filtering
    nan_mask = np.isnan(x)
    if np.any(nan_mask):
        med = np.nanmedian(x)
        x = x.copy()
        x[nan_mask] = med

    # Design notch
    w0 = f0 / (0.5 * fs)
    b, a = iirnotch(w0, Q)
    y = filtfilt(b, a, x)
    # Revert original NaN regions
    y[nan_mask] = np.nan
    return y

ppg_notched = notch_filter(ppg_filled, fs, f0=50.0, Q=30)
abp_notched = notch_filter(abp_filled, fs, f0=50.0, Q=30)
```

**Why notch before baseline removal?**

- 50 Hz (or 60 Hz) spikes can corrupt baseline estimators (e.g., medians). By removing that discrete frequency first, a subsequent spline/median detrend is more stable.

### 2.3.3   Wavelet Denoising (Optional but Powerful)

Wavelet denoising can remove high-frequency noise (e.g., electronic noise, motion spikes remaining after masking). A popular choice is a **db4 discrete wavelet** with soft-thresholding:

```python
import pywt

def wavelet_denoise(x: np.ndarray, wavelet: str="db4", lvl: int=5) -> np.ndarray:
    """
    Perform wavelet denoising on x. Any NaNs are preserved.
    We decompose to 'lvl' levels, threshold detail coefficients, then reconstruct.
    """
    # Preserve NaN mask
    nan_mask = np.isnan(x)
    x0 = x.copy()
    med = np.nanmedian(x0)
    x0[nan_mask] = med  # fill with median

    coeffs = pywt.wavedec(x0, wavelet, level=lvl)
    sigma = np.median(np.abs(coeffs[-1])) / 0.6745  # noise estimate
    uthresh = sigma * np.sqrt(2 * np.log(len(x0)))
    # Soft-threshold detail coefficients
    denoised = [coeffs[0]] + [
        pywt.threshold(c, value=uthresh, mode='soft') for c in coeffs[1:]
    ]
    y = pywt.waverec(denoised, wavelet)
    y = y[: len(x0)]   # truncate if needed
    y[nan_mask] = np.nan
    return y
```

```
ppg_wave = wavelet_denoise(ppg_notched, wavelet="db4", lvl=5)
abp_wave = wavelet_denoise(abp_notched, wavelet="db4", lvl=5)
```

**Why wavelets?**

- Unlike a simple Butterworth low-pass, wavelet denoising can adaptively remove noise at multiple scales.

- Most state-of-the-art PPG pipelines (e.g., Elgendi et al. NPJ Digit. Med. 2019 s41598-023-35492-y (5), and follow-ons) use wavelets to eliminate high-frequency glitches while preserving upstroke/diastolic morphology.

### 2.3.4   Baseline Drift Removal (High-Pass or Median Detrending)

After notch + wavelet, we still see **slow baseline drifts** (respiratory swings, sensor shifts). Two popular, robust approaches:

1. **Butterworth High-Pass** at 0.5 Hz (for PPG) or 0.05 Hz (for ABP).

2. **Median Detrending**: subtract a moving-median over 5 s (for PPG) or 20 s (for ABP).

Below is a combined example:

```python
from scipy.signal import butter, filtfilt

def highpass_filter(x, fs, cutoff, order=4):
    # Fill NaNs temporarily
    nan_mask = np.isnan(x)
    med = np.nanmedian(x)
    x0 = x.copy()
    x0[nan_mask] = med

    nyq = 0.5 * fs
    wn = cutoff / nyq
    b, a = butter(order, wn, btype='highpass')
    y = filtfilt(b, a, x0)
    y[nan_mask] = np.nan
    return y

# For PPG: cutoff = 0.5 Hz
ppg_hp = highpass_filter(ppg_wave, fs, cutoff=0.5, order=4)

# For ABP: cutoff = 0.05 Hz (removes respirations/zeroing drifts)
abp_hp = highpass_filter(abp_wave, fs, cutoff=0.05, order=4)
```

**Alternative median detrending**:

```python
from scipy.signal import medfilt
def median_detrend(x, fs, window_s=5):
    w = int(window_s * fs) | 1  # ensure odd
    nan_mask = np.isnan(x)
    med = np.nanmedian(x)
    x0 = x.copy()
    x0[nan_mask] = med
    baseline = medfilt(x0, kernel_size=w)
    y = x0 - baseline
    y[nan_mask] = np.nan
    return y
```

```
ppg_hp = median_detrend(ppg_wave, fs, window_s=5)
abp_hp = median_detrend(abp_wave, fs, window_s=20)
```

- Median detrending is more robust if there are leftover outliers or calibration spikes that survive wavelet denoising.

Either approach is acceptable; it often depends on which gives the flattest baseline upon visual inspection (compare both on a 1 min stretch to decide). The key is to ensure **no large DC drift** remains.

### 2.3.5   Band-Pass / Low-Pass Finishing Touch

Finally, once the signals are **notched**, **wavelet-denoised**, and **baseline-corrected**, we isolate the physiologic band:

- **PPG band-pass: 0.5–8 Hz** (heart rates between 30–480 bpm, though in ICU typically 40–180 bpm, so an upper cut of 8 Hz is safe).

- **ABP low-pass: < 20 Hz** (ABP wave shape rarely has meaningful content above 20 Hz).

```python
def bandpass_filter(x, fs, lowcut, highcut, order=4):
    nan_mask = np.isnan(x)
    med = np.nanmedian(x)
    x0 = x.copy()
    x0[nan_mask] = med

    nyq = 0.5 * fs
    low = lowcut / nyq
    high = highcut / nyq
    b, a = butter(order, [low, high], btype='bandpass')
    y = filtfilt(b, a, x0)
    y[nan_mask] = np.nan
    return y

def lowpass_filter(x, fs, cutoff, order=4):
    nan_mask = np.isnan(x)
    med = np.nanmedian(x)
    x0 = x.copy()
    x0[nan_mask] = med

    nyq = 0.5 * fs
    wn = cutoff / nyq
    b, a = butter(order, wn, btype='lowpass')
    y = filtfilt(b, a, x0)
    y[nan_mask] = np.nan
    return y

ppg_clean = bandpass_filter(ppg_hp, fs, lowcut=0.5, highcut=8.0, order=4)
abp_clean = lowpass_filter(abp_hp, fs, cutoff=20.0, order=4)
```

**Why zero-phase (`filtfilt`)?**

- It avoids phase distortion of the **systolic upstroke** or the **dicrotic notch**, which downstream feature-extraction (e.g. peak finding) relies on.

At this point, both `ppg_clean` and `abp_clean` should be mostly artifact-free, with a flat baseline, suppressed line noise, and only physiologically plausible oscillations.

## 2.4 Step 4: Quality Assessment & 10 s Segmentation

Having created `ppg_clean` and `abp_clean` at 500 Hz, we now:

1. **Resample (downsample) to 50 Hz** so we can keep consistency with the original CNN design (the PPG2BP-Net used 10 s @ 50 Hz = 500 samples per segment).

2. **Compute beat-quality metrics** (SQIs) on each 10 s window to decide whether to keep it or drop it.

3. **Extract SBP/DBP via robust peak detection** rather than just `max/min`.

### 2.4.1 Downsample to 50 Hz (with Anti-Alias)

```python
from scipy.signal import decimate

def downsample_to_50hz(x: np.ndarray, fs: int) -> np.ndarray:
    """
    Anti-alias decimate from fs (500) -> 50.
    We assume fs % 50 == 0 (i.e., 500 → 50, factor=10).
    Returns a float32 array of length len(x)/decim_factor.
    Any NaNs will be linearly interpolated over short gaps (because decimate cannot handle NaNs),
    or else a fully NaN window stays NaN.
    """
    decim = fs // 50
    # First fill tiny NaN gaps (<0.1s) again
    x_filled = fill_short_gaps(x, fs, max_gap_s=0.1)
    # Now replace remaining NaNs with local median (decimate needs no NaNs)
    nan_mask = np.isnan(x_filled)
    med = np.nanmedian(x_filled)
    x0 = x_filled.copy()
    x0[nan_mask] = med

    y = decimate(x0, decim, ftype="iir", zero_phase=True).astype(np.float32)
    # To preserve artifact regions (that were >0.1s), we can propagate NaNs:
    # map any 10 contiguous input samples that contained any True in nan_mask to one NaN in y
    # e.g. if original index [i*10:(i+1)*10] had any NaN, then y[i] = NaN
    y_nan = np.zeros_like(y, dtype=bool)
    for i in range(len(y)):
        s = i*decim
        e = min(len(x), (i+1)*decim)
        if np.any(nan_mask[s:e]):
            y_nan[i] = True
    y[y_nan] = np.nan
    return y

fs_ds = 50
ppg_ds = downsample_to_50hz(ppg_clean, fs)
abp_ds = downsample_to_50hz(abp_clean, fs)
```

**Key points**:

- We fill only very short NaNs so that decimate sees no gaps.

- We preserve any 200 ms+ artifact as a `NaN` in the output (if any of the 10 input samples was originally a persistent NaN).

- This way, our 50 Hz streams still reveal "holes" where the ABP or PPG was offline/flush-disturbed; we won't let the CNN see those.

### 2.4.2 Divide Into Overlapping 10 s Frames & Compute Quality Metrics

Instead of non-overlapping 10 s windows, it is often beneficial to use a slight **50% overlap** (e.g., stride = 5 s), which increases your dataset and ensures you capture borderline artifacts at window edges. We can then compute an **SQI (Signal Quality Index)** for each 10 s window. Proposed metrics:

1. **Percent NaN**: if more than 5% of samples in a window are NaN → reject.

2. **Coefficient of Variation (CoV)** of the PPG envelope: if CoV < 0.1 → very flat (poor perfusion) → reject.

3. **Beat-Count**: estimate how many beats (via a robust peak-finder) appear in 10 s. If < 5 or > 20 → abnormal heart rates (< 30 bpm or > 120 bpm) → reject.

4. **Baseline Drift Residual**: if the 5 s moving average range > 20 units for PPG or > 50 mmHg for ABP → reject.

Below is a code sketch for sliding-window QC:

```python
from scipy.signal import find_peaks

def compute_window_sqi(ppg_win: np.ndarray, abp_win: np.ndarray, fs_ds: int) -> bool:
    """
    Returns True if this 10s window passes quality checks, False otherwise.
    Checks:
      A) ≤ 5% NaNs in either signal.
      B) PPG envelope CoV between 0.02 and 0.4 (i.e. not too flat, not too noisy).
      C) ABP SBP/DBP in physiologic range (70–180 / 40–110).
      D) Heart rate estimate: 6–18 beats in 10s (36–108 bpm).
      E) Residual baseline drift < threshold.
    """
    Nwin = len(ppg_win)  # should be 10 s * 50 Hz = 500

    # A) NaN fraction
    frac_nan_ppg = np.sum(np.isnan(ppg_win)) / Nwin
    frac_nan_abp = np.sum(np.isnan(abp_win)) / Nwin
    if (frac_nan_ppg > 0.05) or (frac_nan_abp > 0.05):
        return False

    # Fill remaining NaNs with local median for QC metrics
    def fill_for_qc(x):
        nanm = np.isnan(x)
        if np.any(nanm):
            m = np.nanmedian(x)
            x[nanm] = m
        return x

    ppg0 = fill_for_qc(ppg_win.copy())
    abp0 = fill_for_qc(abp_win.copy())

    # B) PPG envelope CoV: take Hilbert transform envelope and check variability
    from scipy.signal import hilbert
    analytic = hilbert(ppg0)
    env = np.abs(analytic)
    CoV = np.std(env) / (np.mean(env) + 1e-6)
    if CoV < 0.02 or CoV > 0.4:  # thresholds roughly tuned
        return False

    # C) ABP SBP/DBP range
    SBP = np.nanmax(abp_win)
    DBP = np.nanmin(abp_win)
    if not (70 <= SBP <= 180 and 40 <= DBP <= 110):
        return False

    # D) Heart rate estimate (ABP peaks)
    peaks, _ = find_peaks(abp0, distance=int(0.3 * fs_ds), height=50)  # at least 0.3 s apart, syst
    num_beats = len(peaks)
    if num_beats < 5 or num_beats > 20:
```

```python
            return False

        # E) Baseline drift residual (5 s median filter)
        w = 5 * fs_ds
        if len(ppg0) >= w:
            baseline_ppg = np.convolve(ppg0, np.ones(w)/w, mode='same')
            drift_ppg = np.nanmax(baseline_ppg) - np.nanmin(baseline_ppg)
            if drift_ppg > 20.0:  # fewer than ±20 units over 5 s is OK; more → reject
                return False
        if len(abp0) >= w:
            baseline_abp = np.convolve(abp0, np.ones(w)/w, mode='same')
            drift_abp = np.nanmax(baseline_abp) - np.nanmin(baseline_abp)
            if drift_abp > 50.0:  # ±50 mmHg drift over 5 s is too much
                return False

        return True

    def sliding_windows_QC(ppg_ds: np.ndarray, abp_ds: np.ndarray, fs_ds: int=50,
                           win_s: float=10.0, stride_s: float=5.0) -> list[tuple[int,int]]:
        """
        Slide a window of win_s (10 s) with stride stride_s (5 s).
        Return a list of (start_idx, end_idx) for each 10 s window that passes QC.
        """
        win_len = int(win_s * fs_ds)        # 500 samples
        stride = int(stride_s * fs_ds)      # 250 samples
        total = len(ppg_ds)
        valid_windows = []

        for start in range(0, total - win_len + 1, stride):
            end = start + win_len
            ppg_win = ppg_ds[start:end]
            abp_win = abp_ds[start:end]
            if compute_window_sqi(ppg_win, abp_win, fs_ds):
                valid_windows.append((start, end))
        return valid_windows

    # Apply:
    valid_idxs = sliding_windows_QC(ppg_ds, abp_ds, fs_ds=50, win_s=10.0, stride_s=5.0)
```

**Why overlapping windows?**

- By stepping every 5 s instead of 10 s, you double your training set and avoid "boundary" frames where a beat is split between two windows (which can destroy SBP/DBP extraction).

- Most subjects in PPG2BP-Net ended up with ~50–100 windows. With overlap, you'll likely collect more valid segments (and can still randomly sub-sample them later).

### 2.4.3 Extract SBP/DBP with Robust Peak Finding

For each valid 10 s window, we can now extract SBP/DBP more accurately:

```python
python

def extract_sbp_dbp(abp_win: np.ndarray, fs_ds: int=50) -> tuple[float,float]:
    """
    Use find_peaks on ABP to locate systolic peaks (height>60 mmHg, distance>0.3 s).
    DBP from troughs (inverted peaks).
    Return mean(SBP_peaks), mean(DBP_troughs). If no peaks/troughs, fallback to max/min.
    """
    abp0 = abp_win.copy()
    # Fill small NaNs before peak finding
    nanm = np.isnan(abp0)
    if np.any(nanm):
        idx = np.arange(len(abp0))
        good = np.where(~nanm)[0]
        if len(good) >= 2:
```

```
            abp0[nanm] = np.interp(idx[nanm], idx[good], abp0[good])
        else:
            # fallback
            return float(np.nanmax(abp_win)), float(np.nanmin(abp_win))

    # 1) Systolic peaks
    peaks, prop = find_peaks(abp0, distance=int(0.3*fs_ds), height=60)
    if len(peaks) >= 3:
        SBP = float(np.mean(abp0[peaks]))
    else:
        SBP = float(np.nanmax(abp0))

    # 2) Diastolic troughs (peak on inverted)
    inv = -abp0
    troughs, prop2 = find_peaks(inv, distance=int(0.3*fs_ds), height=-40)  # only consider < ~40 mmH
    if len(troughs) >= 3:
        DBP = float(np.mean(abp0[troughs]))
    else:
        DBP = float(np.nanmin(abp0))

    return SBP, DBP
```

**Why "≥ 3 peaks/troughs"?**

- In a 10 s window, a normal heart rate of 60–100 bpm gives 10–17 beats. If we find only 1–2 "peaks," it likely means the window is artificially short (maybe it straddles an artifact) or the patient's BP is very low/high. In that case, we fall back to `max`/`min` to get some label, but we also check physiological bounds (T4 below) to decide if this window is valid at all.

### 2.4.4   Window-Level Normalization

Once SBP/DBP are extracted, you can normalize PPG. However, rather than z-scoring each 10 s window independently (which can distort relative amplitudes), you may want to:

1. **Global (per-subject) normalization**: compute μ/σ over *all accepted windows* for that subject and then z-score each window.

2. **Or** keep per-window z-scoring (as in the original paper), but track the scaling factors so that features between windows remain comparable.

Most state-of-the-art pipelines for PPG2BP-like tasks actually do **per-window z-score** (0 mean, 1 SD) because the CNN is supervised to find relative shapes rather than absolute baseline. That is acceptable. Just be careful to compute μ/σ on the **clean (artifact-free) portion** of each 10 s window:

```python

def normalize_ppg_window(ppg_win: np.ndarray) -> np.ndarray:
    """
    Z-score normalize only the non-NaN samples in the 10 s window.
    """
    valid = ~np.isnan(ppg_win)
    if np.sum(valid) < len(ppg_win) * 0.9:  # if >10% is NaN, window should have been rejected alrea
        return np.full_like(ppg_win, np.nan)
    μ = np.nanmean(ppg_win)
    σ = np.nanstd(ppg_win)
    if σ < 1e-6:
        return np.full_like(ppg_win, np.nan)
    out = (ppg_win - μ) / σ
    return out.astype(np.float32)
```

## 2.5   Step 5: Per-Subject Balancing (T5), SDS Calculation, and Splitting

With all "valid windows" (start/end indices) collected, plus their `(normalized_ppg, SBP, DBP)`, you can now:

1. **Group by subject**:

```python
subject_segments = {}  # {caseid: [(ppg_win, SBP, DBP), ...]}
for (start, end) in valid_idxs:
    ppg_win = ppg_ds[start:end]
    abp_win = abp_ds[start:end]
    SBP, DBP = extract_sbp_dbp(abp_win, fs_ds=50)
    # Convert to float and check bounds again (T4)
    if SBP < 70 or SBP > 180 or DBP < 40 or DBP > 110:
        continue

    ppg_norm = normalize_ppg_window(ppg_win)
    if np.isnan(ppg_norm).all():
        continue

    subject_segments[caseid].append((ppg_norm, SBP, DBP))
```

2. **Ensure each subject has ≥ 50 windows**.

   - If `< 50`, drop the subject.

   - If `> 100`, randomly pick 100.

3. **Compute SDS** exactly as your `compute_SDS` function:

```python
from statistics import stdev

def compute_SDS(segments: list[tuple[np.ndarray,float,float]]) -> tuple[float,float]:
    # identical to your implementation
    SBPs = np.array([seg[1] for seg in segments], dtype=np.float32)
    DBPs = np.array([seg[2] for seg in segments], dtype=np.float32)
    if len(SBPs) <= 1:
        return 0.0, 0.0
    SBP_cal, DBP_cal = SBPs[0], DBPs[0]
    delta_SBP = SBPs - SBP_cal
    delta_DBP = DBPs - DBP_cal
    return float(np.std(delta_SBP, ddof=1)), float(np.std(delta_DBP, ddof=1))
```

4. **Subject-Independent 70/10/20 Split** (exactly as before, seed=42).

5. **Save Final** `.npz`: for each subject in train/val/test, save arrays:

```python
np.savez_compressed(
    out_path,
    PPG_segments = np.stack([s[0] for s in segs], axis=0),  # shape (K, 10s*50Hz)
    SBP_labels   = np.array([s[1] for s in segs], dtype=np.float32),
    DBP_labels   = np.array([s[2] for s in segs], dtype=np.float32),
    SBP_cal      = np.float32(segs[0][1]),
    DBP_cal      = np.float32(segs[0][2]),
    SDS_SBP      = np.float32(SDS_SBP),
    SDS_DBP      = np.float32(SDS_DBP),
    age          = np.float32(age),
    sex          = sex,
    weight       = np.float32(weight),
```

```
        height       = np.float32(height)
    )
```

At this point, your `PPG_segments` and `ABP_labels` are derived from **10 s windows that have passed very strict, multi-level QC**. You should end up with even cleaner training data. In particular:

- You have removed all large ABP flushes, zeroing events, and line noise.

- Your PPG no longer contains long flatlines or motion spikes.

- Each 10 s frame has a reliable set of systolic peaks and diastolic troughs so that labels make sense.

- You likely obtain fewer "borderline" windows (fewer false positives), so each sample the CNN sees is truly physiologically plausible.

---

## 3.  Putting It All Together: Pseudocode

Below is a consolidated pseudocode version of the entire **improved** pipeline. You can adapt this to your existing script, replacing the corresponding steps. Wherever you see a comment referencing the PPG2BP-Net paper, the logic tracks the five T-conditions (T1–T5) from Fig. 3 in $^{s}$41598-023-35492-y (5).

```python
#!/usr/bin/env python3
"""
full_preprocess_vitaldb_v2.py

This updated script implements an enhanced PPG2BP-Net preprocessing pipeline,
adding robust artifact detection, notch + wavelet denoising, and strict window-level QC.
- Each subject yields 10s PPG frames @50Hz that have passed multi-stage QC.
- Corresponding SBP/DBP labels are extracted via robust peak finding.
- Output .npz files mirror the original PPG2BP-Net format (50—100 windows per subject).
"""

import os, argparse, random
import numpy as np
import pandas as pd
from scipy.ndimage import label
from scipy.signal import butter, filtfilt, medfilt, decimate, find_peaks
import pywt

# ----------------------------------
# Utility / QC / Denoising Functions
# ----------------------------------

def load_and_qc_subject(caseid, raw_dir, meta_df, min_duration_min=10.0):
    # T1: demographic QC
    row = meta_df[meta_df.caseid == caseid]
    if row.empty:
        return None
    age, wt, ht = float(row.age), float(row.weight), float(row.height)
    if not (18 <= age <= 90 and 10 <= wt <= 100 and 100 <= ht <= 200):
        return None

    # Load raw signals
    path = os.path.join(raw_dir, str(caseid), "signals.npz")
    if not os.path.isfile(path):
        return None
    data = np.load(path)
    ppg_raw = data.get("ppg", None)
    abp_raw = data.get("abp", None)
    fs_raw = float(data.get("fs", 0.0))
    if ppg_raw is None or abp_raw is None or fs_raw != 500.0:
        return None
```

```python
        total_samples = min(len(ppg_raw), len(abp_raw))
        if total_samples < int(min_duration_min * 60 * fs_raw):
            return None

        return ppg_raw[:total_samples].astype(np.float32), abp_raw[:total_samples].astype(np.float32), :

# 2.2.1 ABP Artifact Masking
def mask_abp_artifacts(abp, fs):
    N = len(abp)
    mask = np.ones(N, dtype=bool)

    # (a) Out-of-range
    mask[abp <    0] = False
    mask[abp > 200] = False

    # (b) Flatline detection (std < 1 mmHg over 1s)
    win = fs
    if N >= win:
        from numpy.lib.stride_tricks import sliding_window_view
        windows = sliding_window_view(abp, win)
        stds = np.std(windows, axis=1)
        centers = np.where(stds < 1.0)[0] + win//2
        for c in centers:
            s = max(0, c - win//2)
            e = min(N, c + win//2 + 1)
            mask[s:e] = False

    # (c) Instantaneous spike: Δp > 50 mmHg
    diffs = np.abs(np.diff(abp))
    idx_spike = np.where(diffs > 50.0)[0]
    for i in idx_spike:
        s, e = max(0, i - 10), min(N, i + 10)
        mask[s:e] = False

    return mask

# 2.2.2 PPG Artifact Masking
def mask_ppg_artifacts(ppg, fs):
    N = len(ppg)
    mask = np.ones(N, dtype=bool)

    # (a) Flatline: Δx == 0 for > 100ms
    diffs = np.diff(ppg)
    flat = diffs == 0
    labels_flat, num_flat = label(flat)
    for lab in range(1, num_flat+1):
        idxs = np.where(labels_flat == lab)[0]
        if len(idxs) >= 0.1 * fs:  # >100ms
            s, e = idxs[0], idxs[-1] + 1
            mask[s:e+1] = False

    # (b) Spike: Δ > (μ+3σ)
    abs_d = np.abs(diffs)
    μ, σ = np.mean(abs_d), np.std(abs_d)
    thr = μ + 3*σ
    idx_spike = np.where(abs_d > thr)[0]
    for i in idx_spike:
        s, e = max(0, i - 10), min(N, i + 10)
        mask[s:e] = False

    # (c) Very low amplitude (poor perfusion)
    low_thresh = np.percentile(ppg, 5)
    mask[ppg < low_thresh] = False

    return mask

# 2.3.1 Fill Short Gaps (<0.2s)
def fill_short_gaps(x, fs, max_gap_s=0.2):
    x0 = x.copy().astype(np.float32)
    isnan = np.isnan(x0)
    if not np.any(isnan):
        return x0
    idx = np.arange(len(x0))
    labels_nan, num_nan = label(isnan)
```

```python
    for lab in range(1, num_nan+1):
        idxs = np.where(labels_nan == lab)[0]
        if len(idxs) <= max_gap_s * fs:
            good = np.where(~isnan)[0]
            x0[idxs] = np.interp(idxs, good, x0[good])
    return x0

# 2.3.2 Notch Filter (50Hz)
def notch_filter(x, fs, f0=50.0, Q=30.0):
    x0 = x.copy().astype(np.float32)
    nanm = np.isnan(x0)
    if np.any(nanm):
        med = np.nanmedian(x0)
        x0[nanm] = med

    w0 = f0 / (0.5 * fs)
    b, a = butter(2, w0, btype='bandstop')  # second-order notch
    y = filtfilt(b, a, x0)
    y[nanm] = np.nan
    return y

# 2.3.3 Wavelet Denoise
def wavelet_denoise(x, wavelet='db4', lvl=5):
    nanm = np.isnan(x)
    x0 = x.copy().astype(np.float32)
    if np.any(nanm):
        med = np.nanmedian(x0)
        x0[nanm] = med

    coeffs = pywt.wavedec(x0, wavelet, level=lvl)
    sigma = np.median(np.abs(coeffs[-1])) / 0.6745
    uthresh = sigma * np.sqrt(2 * np.log(len(x0)))
    coeffs_denoised = [coeffs[0]] + [pywt.threshold(c, uthresh, mode='soft') for c in coeffs[1:]]
    y = pywt.waverec(coeffs_denoised, wavelet)[:len(x0)]
    y[nanm] = np.nan
    return y

# 2.3.4 High-Pass Filter (Butterworth)
def highpass_filter(x, fs, cutoff, order=4):
    x0 = x.copy().astype(np.float32)
    nanm = np.isnan(x0)
    if np.any(nanm):
        med = np.nanmedian(x0)
        x0[nanm] = med

    nyq = 0.5 * fs
    wn = cutoff / nyq
    b, a = butter(order, wn, btype='highpass')
    y = filtfilt(b, a, x0)
    y[nanm] = np.nan
    return y

# 2.3.5 Bandpass & Lowpass
def bandpass_filter(x, fs, lowcut, highcut, order=4):
    x0 = x.copy().astype(np.float32)
    nanm = np.isnan(x0)
    if np.any(nanm):
        med = np.nanmedian(x0)
        x0[nanm] = med

    nyq = 0.5 * fs
    low = lowcut / nyq
    high = highcut / nyq
    b, a = butter(order, [low, high], btype='bandpass')
    y = filtfilt(b, a, x0)
    y[nanm] = np.nan
    return y

def lowpass_filter(x, fs, cutoff, order=4):
    x0 = x.copy().astype(np.float32)
    nanm = np.isnan(x0)
    if np.any(nanm):
        med = np.nanmedian(x0)
        x0[nanm] = med
```

```python
    nyq = 0.5 * fs
    wn = cutoff / nyq
    b, a = butter(order, wn, btype='lowpass')
    y = filtfilt(b, a, x0)
    y[nanm] = np.nan
    return y

# 2.4.1 Downsample to 50Hz
def downsample_to_50hz(x, fs):
    decim = fs // 50
    x_filled = fill_short_gaps(x, fs, max_gap_s=0.1)
    nanm = np.isnan(x_filled)
    if np.any(nanm):
        med = np.nanmedian(x_filled)
        x_filled[nanm] = med

    y = decimate(x_filled, decim, ftype='iir', zero_phase=True).astype(np.float32)
    # Propagate NaNs: if any of the decim input block was NaN, we set the output sample to NaN
    y_nan = np.zeros_like(y, dtype=bool)
    for i in range(len(y)):
        s = i * decim
        e = min(len(x), (i+1) * decim)
        if np.any(nanm[s:e]):
            y_nan[i] = True
    y[y_nan] = np.nan
    return y

# 2.4.2 Window-Level QC
from scipy.signal import hilbert

def compute_window_sqi(ppg_win, abp_win, fs_ds):
    Nwin = len(ppg_win)
    # A) NaN fraction
    if np.sum(np.isnan(ppg_win)) / Nwin > 0.05 or np.sum(np.isnan(abp_win)) / Nwin > 0.05:
        return False

    # Fill NaNs with local median
    def fill_for_qc(xx):
        mxx = xx.copy()
        nanm = np.isnan(mxx)
        if np.any(nanm):
            med = np.nanmedian(mxx)
            mxx[nanm] = med
        return mxx

    ppg0 = fill_for_qc(ppg_win.copy())
    abp0 = fill_for_qc(abp_win.copy())

    # B) PPG envelope CoV
    analytic = hilbert(ppg0)
    env = np.abs(analytic)
    CoV = np.std(env) / (np.mean(env) + 1e-6)
    if CoV < 0.02 or CoV > 0.4:
        return False

    # C) ABP SBP/DBP in range
    SBP = float(np.nanmax(abp_win))
    DBP = float(np.nanmin(abp_win))
    if not (70 <= SBP <= 180 and 40 <= DBP <= 110):
        return False

    # D) Heart rate estimate via ABP peaks
    peaks, _ = find_peaks(abp0, distance=int(0.3 * fs_ds), height=50)
    if len(peaks) < 5 or len(peaks) > 20:
        return False

    # E) Residual baseline drift
    w = 5 * fs_ds
    if len(ppg0) >= w:
        baseline_ppg = np.convolve(ppg0, np.ones(w)/w, mode='same')
        if (np.nanmax(baseline_ppg) - np.nanmin(baseline_ppg)) > 20.0:
            return False
    if len(abp0) >= w:
        baseline_abp = np.convolve(abp0, np.ones(w)/w, mode='same')
        if (np.nanmax(baseline_abp) - np.nanmin(baseline_abp)) > 50.0:
```

```python
            return False

    return True

def sliding_windows_QC(ppg_ds, abp_ds, fs_ds=50, win_s=10.0, stride_s=5.0):
    win_len = int(win_s * fs_ds)
    stride  = int(stride_s * fs_ds)
    total   = len(ppg_ds)
    valid   = []
    for start in range(0, total - win_len + 1, stride):
        end = start + win_len
        if compute_window_sqi(ppg_ds[start:end], abp_ds[start:end], fs_ds):
            valid.append((start, end))
    return valid

# 2.4.3 Robust SBP/DBP Extraction
def extract_sbp_dbp(abp_win, fs_ds=50):
    """
    Return (SBP, DBP). Use actual peaks if >=3 found, otherwise fallback to max/min.
    """
    aw = abp_win.copy()
    nanm = np.isnan(aw)
    if np.any(nanm):
        idx = np.arange(len(aw))
        good = np.where(~nanm)[0]
        if len(good) >= 2:
            aw[nanm] = np.interp(idx[nanm], idx[good], aw[good])
        else:
            return float(np.nanmax(abp_win)), float(np.nanmin(abp_win))

    peaks, _ = find_peaks(aw, distance=int(0.3*fs_ds), height=60)
    if len(peaks) >= 3:
        SBP = float(np.mean(aw[peaks]))
    else:
        SBP = float(np.nanmax(aw))

    troughs, _ = find_peaks(-aw, distance=int(0.3*fs_ds), height=-40)
    if len(troughs) >= 3:
        DBP = float(np.mean(aw[troughs]))
    else:
        DBP = float(np.nanmin(aw))

    return SBP, DBP

# 2.4.4 Normalize PPG Window
def normalize_ppg_window(ppg_win):
    valid = ~np.isnan(ppg_win)
    if np.sum(valid) < len(ppg_win) * 0.9:
        return np.full_like(ppg_win, np.nan)
    μ = float(np.nanmean(ppg_win))
    σ = float(np.nanstd(ppg_win))
    if σ < 1e-6:
        return np.full_like(ppg_win, np.nan)
    return ((ppg_win - μ) / σ).astype(np.float32)

# 2.5 SDS Calculation
def compute_SDS(segments):
    SBPs = np.array([s[1] for s in segments], dtype=np.float32)
    DBPs = np.array([s[2] for s in segments], dtype=np.float32)
    if len(SBPs) <= 1:
        return 0.0, 0.0
    SBP_cal, DBP_cal = SBPs[0], DBPs[0]
    delta_SBP = SBPs - SBP_cal
    delta_DBP = DBPs - DBP_cal
    return float(np.std(delta_SBP, ddof=1)), float(np.std(delta_DBP, ddof=1))


# -------------------------------
# Main Preprocessing Routine
# -------------------------------

def full_preprocess_v2(raw_dir, meta_csv, out_dir, min_duration_min=10.0, fs_target=50):
    """
    Implements T1–T5 with advanced artifact handling and QC.
    """
```

```python
np.random.seed(42)
random.seed(42)

# 1) Load metadata & apply T1
meta_df = pd.read_csv(meta_csv)
meta_df = meta_df[(meta_df.age.between(18,90)) &
                  (meta_df.weight.between(10,100)) &
                  (meta_df.height.between(100,200))].copy()
meta_df.caseid = meta_df.caseid.astype(int)

balanced_segments = {}  # { caseid: [(ppg_norm, SBP, DBP), ...] }
sds_dict = {}           # { caseid: (SDS_SBP, SDS_DBP) }
dropped_t2 = dropped_t3 = dropped_t5 = 0

all_caseids = meta_df.caseid.tolist()
for caseid in all_caseids:
    # T1/T2: load & QC
    tmp = load_and_qc_subject(caseid, raw_dir, meta_df, min_duration_min)
    if tmp is None:
        dropped_t2 += 1
        continue
    ppg_raw, abp_raw, fs_raw, age, sex, wt, ht = tmp

    # 2) Mask artifacts (PPG & ABP)
    ppg_mask = mask_ppg_artifacts(ppg_raw, fs_raw)
    abp_mask = mask_abp_artifacts(abp_raw, fs_raw)
    combined_mask = ppg_mask & abp_mask

    ppg_masked = ppg_raw.copy().astype(np.float32);  ppg_masked[~combined_mask] = np.nan
    abp_masked = abp_raw.copy().astype(np.float32);  abp_masked[~combined_mask] = np.nan

    # 3) Fill short gaps, notch, wavelet, HPF, bandpass/lowpass
    ppg_filled = fill_short_gaps(ppg_masked, fs_raw, max_gap_s=0.2)
    abp_filled = fill_short_gaps(abp_masked, fs_raw, max_gap_s=0.2)

    ppg_notch = notch_filter(ppg_filled, fs_raw, f0=50.0, Q=30)
    abp_notch = notch_filter(abp_filled, fs_raw, f0=50.0, Q=30)

    ppg_wave = wavelet_denoise(ppg_notch, wavelet='db4', lvl=5)
    abp_wave = wavelet_denoise(abp_notch, wavelet='db4', lvl=5)

    ppg_hp = highpass_filter(ppg_wave, fs_raw, cutoff=0.5, order=4)
    abp_hp = highpass_filter(abp_wave, fs_raw, cutoff=0.05, order=4)

    ppg_clean = bandpass_filter(ppg_hp, fs_raw, lowcut=0.5, highcut=8.0, order=4)
    abp_clean = lowpass_filter(abp_hp, fs_raw, cutoff=20.0, order=4)

    # 4) Downsample to 50Hz
    ppg_ds = downsample_to_50hz(ppg_clean, fs_raw)
    abp_ds = downsample_to_50hz(abp_clean, fs_raw)
    fs_ds = fs_target  # now 50 Hz

    # 5) T3+T4+window-level QC & Segment
    valid_windows = sliding_windows_QC(ppg_ds, abp_ds, fs_ds, win_s=10.0, stride_s=5.0)
    if not valid_windows:
        dropped_t3 += 1
        continue

    # Build a list of (ppg_norm, SBP, DBP)
    segs = []
    for (s, e) in valid_windows:
        pwin = ppg_ds[s:e]
        awin = abp_ds[s:e]

        SBP, DBP = extract_sbp_dbp(awin, fs_ds)
        if (SBP < 70 or SBP > 180) or (DBP < 40 or DBP > 110):
            continue

        ppg_norm = normalize_ppg_window(pwin)
        if np.isnan(ppg_norm).all():
            continue

        segs.append((ppg_norm, SBP, DBP))

    if len(segs) < 50:
```

```python
            dropped_t5 += 1
            continue

        # 2.5.1 T5: balance segments (50–100)
        if len(segs) > 100:
            segs = random.sample(segs, 100)

        balanced_segments[caseid] = segs
        sds_dict[caseid] = compute_SDS(segs)

print("== Preprocessing Summary ==")
print(f"Subjects dropped at T2: {dropped_t2}")
print(f"Subjects dropped at T3 (no valid windows): {dropped_t3}")
print(f"Subjects dropped at T5 (<50 windows): {dropped_t5}")
print(f"Subjects remaining: {len(balanced_segments)}")

# 6) Build DataFrame of surviving subjects
rows = []
for caseid, segs in balanced_segments.items():
    r = meta_df[meta_df.caseid == caseid].iloc[0]
    SDS_SBP, SDS_DBP = sds_dict[caseid]
    rows.append({
        "caseid": caseid,
        "age": float(r.age),
        "sex": r.sex,
        "weight": float(r.weight),
        "height": float(r.height),
        "num_segments": len(segs),
        "SDS_SBP": SDS_SBP,
        "SDS_DBP": SDS_DBP
    })
df_all = pd.DataFrame(rows).sort_values("caseid").reset_index(drop=True)

# 7) 70/10/20 Subject-Independent Split
all_cids = df_all.caseid.tolist()
random.shuffle(all_cids)
n_total = len(all_cids)
n_train = int(0.7 * n_total)
n_val   = int(0.1 * n_total)
train_cids = all_cids[:n_train]
val_cids   = all_cids[n_train:n_train+n_val]
test_cids  = all_cids[n_train+n_val:]

# 8) Save per-subject .npz
for split, cids in [("train", train_cids), ("val", val_cids), ("test", test_cids)]:
    split_dir = os.path.join(out_dir, split)
    os.makedirs(split_dir, exist_ok=True)
    for cid in cids:
        segs = balanced_segments[cid]
        PPGs = np.stack([s[0] for s in segs], axis=0)  # (K, 500)
        SBPs = np.array([s[1] for s in segs], dtype=np.float32)
        DBPs = np.array([s[2] for s in segs], dtype=np.float32)
        SDS_SBP, SDS_DBP = sds_dict[cid]

        r = meta_df[meta_df.caseid == cid].iloc[0]
        SBP_cal = segs[0][1]
        DBP_cal = segs[0][2]

        np.savez_compressed(
            os.path.join(split_dir, f"{cid}.npz"),
            PPG_segments = PPGs,
            SBP_labels   = SBPs,
            DBP_labels   = DBPs,
            SBP_cal      = np.float32(SBP_cal),
            DBP_cal      = np.float32(DBP_cal),
            SDS_SBP      = np.float32(SDS_SBP),
            SDS_DBP      = np.float32(SDS_DBP),
            age          = np.float32(r.age),
            sex          = r.sex,
            weight       = np.float32(r.weight),
            height       = np.float32(r.height)
        )

# 9) Save CSV summaries
os.makedirs(out_dir, exist_ok=True)
```

```
        df_all.to_csv(os.path.join(out_dir, "all_subjects_info.csv"), index=False)
        pd.DataFrame({"train_cid": train_cids}).to_csv(os.path.join(out_dir, "train_cids.csv"), index=Fa
        pd.DataFrame({"val_cid":   val_cids}).to_csv(os.path.join(out_dir, "val_cids.csv"), index=False
        pd.DataFrame({"test_cid":  test_cids}).to_csv(os.path.join(out_dir, "test_cids.csv"), index=Fals

        print("=== Finished enhanced preprocessing. Outputdir:", out_dir, "===")
        print(f"Subjects: train={len(train_cids)}, val={len(val_cids)}, test={len(test_cids)}")
        print("Avg SDS_SBP (train) =", df_all[df_all.caseid.isin(train_cids)]["SDS_SBP"].mean())
        print("Avg SDS_SBP (val)   =", df_all[df_all.caseid.isin(val_cids)]["SDS_SBP"].mean())
        print("Avg SDS_SBP (test)  =", df_all[df_all.caseid.isin(test_cids)]["SDS_SBP"].mean())


if __name__ == "__main__":
    parser = argparse.ArgumentParser(description="Enhanced PPG2BP-Net preprocessing (T1–T5 + advance
    parser.add_argument("--raw_dir",   required=True,
                        help="Root folder containing raw_data/<caseid>/signals.npz")
    parser.add_argument("--meta_csv", required=True,
                        help="Path to metadata CSV (with columns: caseid, age, sex, weight, height,
    parser.add_argument("--out_dir",   required=True,
                        help="Output folder; will contain train/, val/, test/ subfolders.")
    parser.add_argument("--min_duration_min", type=float, default=10.0,
                        help="Minimum recording length (minutes) to keep a subject (default=10).")
    args = parser.parse_args()

    full_preprocess_v2(
        raw_dir=args.raw_dir,
        meta_csv=args.meta_csv,
        out_dir=args.out_dir,
        min_duration_min=args.min_duration_min,
        fs_target=50
    )
```

---

## 4.   Why This Pipeline Is "State-of-the-Art"

1. **Early Artifact Masking**
   - We explicitly detect ABP flushes, zeroing events, and PPG motion spikes *before* any smoothing.
   - In the original PPG2BP-Net paper, segment-level QC (T3/T4) occurred only *after* downsampling and low-pass. By bringing artifact detection forward, we ensure *filters do not smear* artifact energy into adjacent samples.

2. **Notch + Wavelet Denoising**
   - A simple 25 Hz low-pass cannot remove 50/60 Hz mains noise. Our notch filter at 50 Hz (Q=30) removes power-line interference.
   - Wavelet denoising (db4, 5-level) adaptively removes high-frequency transients (electro-surgical spikes, probe taps) while preserving morphological features like the dicrotic notch. This is a common "best practice" in PPG literature (e.g., Elgendi et al. in NPJ Digit. Med. 2019).

3. **Robust Baseline Removal**
   - We offer both Butterworth high-pass (0.5 Hz for PPG, 0.05 Hz for ABP) *and* median detrending alternatives. Median filters are more robust to outliers than polynomial detrending.

4. **Strict Window-Level QC (SQI)**
   - Instead of simply rejecting windows with any NaN or extreme SBP/DBP, we compute a multi-metric SQI:
     - Fraction of NaNs < 5% (any > 5% = reject).
     - PPG envelope CoV within a physiologic range (0.02–0.4).

- ABP SBP/DBP in range [70–180]/[40–110].

- Reasonable heartbeat count (6–20 beats in 10 s).

- Residual baseline < threshold.

- This dramatically reduces the chance that a "noisy" window slips through.

5. **Overlap & Data Augmentation**

- By sliding windows every 5 s (rather than non-overlap), we (a) double our dataset size, (b) avoid boundary effects, and (c) give the CNN more slightly-shifted examples for robustness.

6. **Robust SBP/DBP Labeling**

- We look for ≥ 3 peaks/troughs within a 10 s window to compute the *average* SBP/DBP, instead of simply `max/min`. This mitigates a single extreme spike (say, from a brief motion) skewing your label.

7. **Better NaN Handling**

- We only fill very short gaps (< 0.1 s or < 0.2 s). Anything longer remains NaN and is ultimately rejected by the QC step. This preserves real "holes" in the data, preventing filters from smoothing them into fake data.

8. **Transparent SDS & Splitting**

- We compute SDS exactly per the paper's definition, ensuring our train/val/test splits replicate the same intrasubject variability criteria.

---

# 5.  Next Steps & Tips

1. **Visual Inspection**

- After implementing the above pipeline, plot a few random 1 min traces of both raw vs. cleaned PPG/ABP to confirm artifacts have been removed but physiological morphology is intact.

- Overlap the raw and the final 50 Hz traces (e.g., `plot(t, raw[fs*60:fs*61], label="raw"); plot(t_ds, ds[60*50:61*50], label="cleaned")`) to visually confirm.

2. **Tune Thresholds**

- The CoV thresholds (0.02–0.4) and drift thresholds (±20 units for PPG, ±50 mm Hg for ABP) are "reasonable defaults." Feel free to tweak slightly depending on your specific sensor/ICU environment.

- If you see too many windows being rejected globally, you may need to loosen the drift thresholds or widen the HR range.

3. **Optional Further Denoising**

- You could add an **Empirical Mode Decomposition (EMD)** step (e.g., EEMD) to separate slow trends from pulsatile components, but wavelet denoising is usually sufficient.

- Another advanced method is to use a small **denoising autoencoder** on PPG/ABP. However, that becomes circular (you need clean data to train it).

4. **Augmentation for Training**

- Once you have your clean 10 s PPG windows, you can consider data augmentation (adding small gaussian noise, random time warping, random cropping to 9 s + padding) to improve CNN generalization, especially if your dataset is imbalanced across SDS levels.

5. **Evaluate Final SDS/Distribution**

- Compute the final SDS distribution for train/val/test. Ideally, they should all cover a similar range (e.g., mean SDS_SBP ≈ 19 mm Hg for train, ≈ 19 mm Hg for val/test) to ensure comparability.

---

## References

- Joung *et al.* (2023) "Continuous cuffless blood pressure monitoring using photoplethysmography-based PPG2BP-Net for high intrasubject blood pressure variations" s41598-023-35492-y (5)

- Elgendi *et al.* (2019) "The use of photoplethysmography for assessing hypertension," NPJ Digit. Med. 2, 1–11.

- Mousavi *et al.* (2019) "Blood pressure estimation from appropriate and inappropriate PPG signals using a whole-based method," Biomed. Signal Process. Control 47, 196–206.

- Liang *et al.* (2018) "Photoplethysmography signal reconstruction by discrete wavelet transform," IEEE Trans. Biomed. Eng. (example of wavelet usage).

---

**In summary**, this enhanced preprocessing pipeline:

1. Masks out large ABP flushes & PPG motion artifacts early.

2. Uses a notch filter + wavelet denoising to suppress line noise & transients.

3. Removes baseline drift with either Butterworth HPF or median detrending.

4. Applies a strict multi-metric window QC to ensure only physiologically valid 10 s frames survive.

5. Falls back to robust peak-finding for SBP/DBP label extraction.

Implementing this will give your CNN the **cleanest possible PPG & ABP inputs**, maximizing its ability to learn the subtle pulse-shape changes needed to estimate blood pressure accurately—especially in high-SDS (high intrasubject variability) conditions.