

Step 1 - which kind of metric?

- Full-ref** Have the clean original? Use a full-reference metric (PESQ, POLQA, ViSQOL, PEAQ). Best in the build pipeline and lab.
- No-ref** No original (live traffic)? Use a no-reference metric (DNSMOS P.835, NISQA). Built from machine learning.
- Truth** All of them predict a Mean Opinion Score (MOS). The room full of human listeners is still the ground truth.

Step 2 - pick the full-reference metric

Metric	Best for	Range	2026 status	Cost
PESQ (P.862)	speech / telephony	to 7 kHz	deleted Jan 2024	free
POLQA (P.863)	speech / telephony	to 20 kHz	current	licensed
ViSQOL v3	speech + music	16 / 48 kHz	maintained	free (Apache 2.0)
PEAQ (BS.1387)	music codecs	full-band	current (2023)	standard

Key facts that catch people

- PESQ was withdrawn in 2018 and deleted from the ITU catalogue on 5 Jan 2024 - yet ~4,600 papers still used it in 2024 (free; POLQA is licensed).
- POLQA fixes PESQ's two big gaps: it reaches full-band audio and scores time-warping (codec/jitter speed changes) the way a listener hears it.
- ViSQOL turns audio into a spectrogram and compares the pictures with NSIM (an image-similarity measure). Audio mode 48 kHz is the rare free music-capable option.
- PEAQ outputs an Objective Difference Grade (ODG), 0 (imperceptible) to -4 (very annoying), for codecs like AAC/MP3.
- 'AAC-Q' is not a real standard - it means music-codec quality, done by PEAQ / ViSQOLAudio / vendor tools.

Using a metric in CI - the discipline

- Match the metric to your content - speech metric on music is confident nonsense.
- Fix a 30-50 clip test set the model never sees; store per-clip scores, not just the mean.
- Set the CI gate from statistics: fail only on a drop bigger than the standard error (s / \sqrt{N}).
- Never train a model on the metric - Goodhart's law makes the score gameable.
- Back it with a periodic P.808 human listening test; if metric and humans disagree, humans win.