

The three regimes, the arithmetic, and the speaker-cap rule - on one page.

Three regimes: 50 / 500 / 5,000 at a glance

PROPERTY	<= 50	~500	5,000+
Architecture	SFU, forward all	SFU, forward active	MCU mix few; cascade
Server decodes audio?	Never	Never	Only active few
Each listener receives	Every voice	Active speakers	One mixed stream
Dominant cost	Client decode	Server fan-out	Server mix CPU
Streams downloaded	1 per talker	1 per talker	Exactly 1 (mix)
Added latency	Lowest	Low	Higher (mix)
Per-speaker control	Full	Full	Lost once mixed
Single-file recording	Separate mix	Separate mix	Mix already exists

Why it explodes: $K \times (K-1)$ voice deliveries

5 people	$5 \times 4 = 20$ voice deliveries.
50 people	$50 \times 49 = 2,450$ - about 122x the 5-person load.
500 people	$500 \times 499 = 249,500$ deliveries if you forward everything.
5,000 people	$5,000 \times 4,999 = 24,995,000$ - quadratic growth: doubling ~4x the work.

The rules that make scale possible

Audio != video	Last-N forwards only top videos; ALL audio is forwarded in small/medium rooms.
DTX	A silent Opus mic drops from ~24-32 kbit/s to ~1 kbit/s - cost follows talkers, not total.
RFC 6464	Loudness stamped per packet (0-127 -dBov); server reads it without decoding.
N-1 mix	Each speaker hears everyone EXCEPT themselves; mix only the loudest 3-6 channels.
Cascade / hybrid	Beyond one server: relay across a mesh, or stream the mix over LL-HLS to passives.

Remember

- <= 50 people: SFU forwards all audio to everyone; DTX keeps muted mics nearly free.
- 50-500 people: forward only the active speakers, ranked by RFC 6464 loudness levels.
- 5,000+ people: mix the loudest 3-6 speakers into ONE stream; never sum all channels.
- Beyond one server's ceiling: cascade across a mesh, or go hybrid (mix + LL-HLS to passives).
- Size for the PEAK simultaneous-talker count, not the average - the 'any questions?' spike.
- Cap simultaneously forwarded/mixed speakers at 3-6; back it with a raise-hand or stage UI.
- Active-speaker UI flickers when people pause? Suspect Opus DTX gaps, not the algorithm.
- If almost nobody needs to talk back, it's a broadcast - serve the passive tier over streaming.