

# JACKTRIP: UNDER THE HOOD OF AN ENGINE FOR NETWORK AUDIO

*Juan-Pablo Cáceres & Chris Chafe*

Center for Computer Research in Music and Acoustics (CCRMA)

Stanford University

{[jcaceres,cc](mailto:jcaceres,cc@ccrma.stanford.edu)}@ccrma.stanford.edu

## ABSTRACT

The design of a platform for bi-directional musical performance using modern WAN networks poses several challenges that are different from related applications, e.g., synchronous LAN studio systems or uni-directional WAN streaming. The need to minimize as much as possible audio latency and also maximize audio quality requires specific strategies which are informed, in part, by musical decisions.

We present some of the key design elements of the JackTrip application which has evolved through several years of deployment in musical work over wide-area networks.

## 1. INTRODUCTION

The SoundWIRE group at CCRMA<sup>1</sup> focuses on experiments with bi-directional musical performance. Concerts and rehearsals between Stanford and places like New York, Belfast, Banff, Beijing, or Santiago are now routine.

JackTrip is the application which powers up these on-line collaborations. Presently, it's a Linux and Mac OS X-based system which supports multi-machine network performance over best-effort Internet. The technology being used builds on early work by research groups at McGill University [11] and Stanford University [7]. The basic approach is to send uncompressed audio (avoiding the latency introduced by compression encode/decode algorithms) through high-speed links like *Internet2*. It supports any number of channels (as many as the computers or network paths can handle). Since best-effort network protocols are used, adequate network provisioning is a must.

The subject of this article is JackTrip's design relating to several issues that come up in implementing such a system. It is hoped that these solutions can serve as a point-of-departure for further applications in this same area.

The design achieves (i) the highest audio quality possible, by using uncompressed linear sampling and redundancy to recover from packet loss; (ii) throughput maximization, which gets audio packets onto and off of the network as soon as the sound card can deliver them; (iii) working with any

number of channels (depending on available computer processing power and bandwidth); (iv) flexibility in routing and mixing audio channels from and to the different hosts.

### 1.1. Peer-to-peer Network Audio Latency

WAN connections inevitably introduce transmission delays between two or more hosts. For non-interactive and "soft" real-time applications, this delay is less of a problem than for high-quality collaborative music performance. The latter places extremely stringent bounds on latency and jitter. The longer the audio latency between musicians, the harder it is for them to play synchronously [5]. Time delays as short as 25 milliseconds are already problematic for professional ensembles like string quartets.<sup>2</sup>

It's the total delay between sound capture and sound projection which counts. This splits out into (i) acoustic (air path) delays, e.g., the distance between an instrument and the capture microphone and between the speakers and ears; (ii) analog to digital and digital to analog conversion (ADC/DAC) delay, i.e., the time it takes for an analog source to be transformed into digital and back; (iii) settings chosen for audio quality and packetization, including audio sampling rate and bit depth resolution, buffer and packet sizes, and others; (iv) network transmission delays, including physical (geographical) distance, transmission delays induced by switches, routers, firewall and network congestion among others.

The default transport protocol in JackTrip is UDP, a low-overhead, fast mechanism for transmitting packets (see [9] for a good description). The application's own header data accompanies each audio packet to describe local properties like audio buffer size, sampling rate, bit depth, number of channels, a sequence number and a time stamp.

Currently, JackTrip uses Jack [3] as its host audio server. Jack has several advantages: it runs on Linux and Mac OS X, it has the ability to make audio connections between many different audio clients on the same host, and its current implementation takes advantage of multi-processor machines [10].

<sup>1</sup><http://ccrma.stanford.edu/groups/soundwire/>.

<sup>2</sup>Recordings of experiments with the *St. Lawrence String Quartet* are available at <http://ccrma.stanford.edu/groups/soundwire/research/slsq/>.

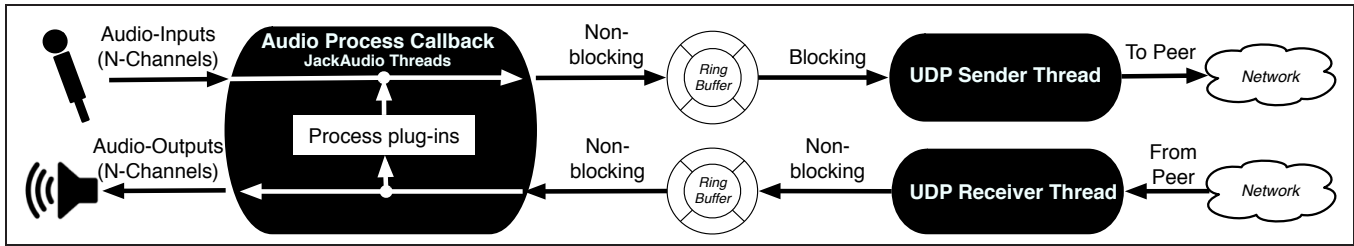


Figure 1. JackTrip architecture overview

## 2. JACKTRIP'S MULTI-THREADED ARCHITECTURE

JackTrip's multi-threaded design is implemented in C++ using the Qt library<sup>3</sup> and also can take advantage of multi-core machines. Figure 1 shows the multi-threaded architecture of the application. There is one thread which processes Jack's audio via a callback function. The other processing threads in JackTrip are the *Sender*, which wraps audio packets from the audio thread with the header information into UDP packets, and the *Receiver* that unwraps the packet and has it ready when the audio process callback needs it.

Inter-thread communication is implemented using ring (or circular) buffers as shown in Figure 1. This is one of the critical latency-reducing parts of the design. The only thread which blocks against its input from the ring buffer is the UDP Sender thread; there's no need to send audio that hasn't been generated. Every time a buffer is available on the ring buffer, the sender thread immediately sends it as a UDP packet. Conversely, the receiving ring buffer cannot block, since local audio must obtain a packet from the ring buffer when it's scheduled to—otherwise audio glitches will be heard. JackTrip maximizes reliability, audio flexibility and minimizes as much as possible peer-to-peer audio delay. Two parameters which affect local audio latency are sampling rate and buffer size. For example, using a sampling rate of 96 kHz and an audio buffer of 64 frames (or samples), the rate of audio packet delivery is every  $64/96000 \cdot 1000 = 0.67$  milliseconds<sup>4</sup>.

### 2.1. Thread Scheduling

Threads in JackTrip are scheduled as *real-time* priority, i.e., jack audio and socket threads will take priority over any other non-critical process. This avoids interruptions during time-critical tasks.

<sup>3</sup><http://www.qtsoftware.com/>.

<sup>4</sup>Internal redundancy and other factors can make the actual local latency approximately the double of this number, but the delivery rate, i.e., the rate at which packets are sent and received, corresponds to that number.

### 2.2. Buffering Mechanism

Two types of scheduling problems can occur on the receive side, illustrated in Figure 2:

**Overrun condition** The receiving ring buffer is full, i.e., there is no space to write new buffers coming from the UDP Receiver thread. This normally happens when asynchronous clocks drift, e.g., the peer's clock runs faster than the local clock.

**Under-run condition** The receiving ring buffer is empty, i.e., there are no new packets to read. This is caused either because there are packets that are delayed or lost in the network or because the clocks of the two machines have drifted the other way.

This is different from common streaming applications which can stop playback (e.g., audio-video playback on browsers) when they reach an under-run condition, and won't have the overflow problem because real-time is not a concern. Typically, these applications adaptively increase or reduce their buffer size. In JackTrip, latency needs to be constant and another method is needed to deal with both under and over-run conditions.

Ring buffers have a *read* and a *write* pointer (Fig. 2). On initialization, both the read and write pointers are in a "symmetric" position. The longer the buffer size, the higher the latency<sup>5</sup>. The length of the buffer is a provision for network jitter; slight variations around this symmetry are produced by packets not arriving at exactly the same frequency. In an ideal situation, where both machines have clocks that match exactly and no packets are lost, the symmetric position will be maintained on average throughout the connection. The higher the jitter the longer the ring buffer needs to be to avoid glitches.

### 2.3. Buffer Glitches

Primarily as a result of receive buffers not being sized to accommodate network jitter glitches occur and have to be dealt with. The application has two different modes that respond to under-runs (Fig. 2):

<sup>5</sup>This latency can be visualized as the signed "distance" between the *read* and *write* pointer.

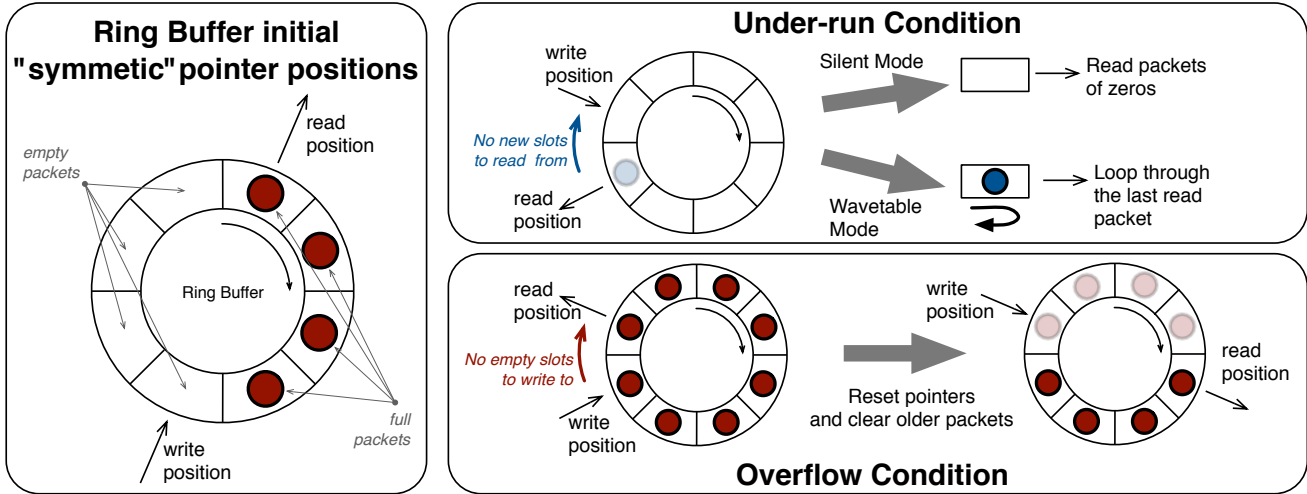


Figure 2. Ring Buffers

**Silent mode** Send a packet of zeros (silence) to the process callback.

**Wavetable mode** Re-send the last available packet to the process callback. This will produce a wavetable synthesizer type of sound when there aren't new packets available for some time, since it's going to loop on the last one received.

For under-runs, the pointers are not reset because we always want to be able to *read* the most recent packet.

To deal with buffer overflows, the ring buffer *read pointer* is reset to the symmetric position with respect to the *write pointer*. Some packets<sup>6</sup> will be lost in the process but the clock drift will be reset to its original position, thus avoiding another glitch for an extended amount of time.

## 2.4. Packet Redundancy Algorithm

As an unreliable transport mechanism, UDP has no provisions to notify the sender when or if a packet was successfully delivered, or if the receiving order matches the sender's. With today's good network QoS, we generally experience a very low number of lost or out-of-order packets. But, since even one misplaced packet will be perceived as a glitch, JackTrip includes a mechanism to recover (within certain bounds) lost or unordered packets.

Jacktrip's redundancy algorithm is used when sufficient bandwidth is available. The technique is illustrated on Figure 3. The sender bundles *RedunFactor* copies of every audio+header packet into a bigger UDP packet (with  $RedunFactor \in \mathbb{Z}^+$ ). This is done for every new audio+header buffer, so each UDP packet has an overlap of  $RedunFactor - 1$  buffers, as illustrated in the figure.

<sup>6</sup>Half of the ring buffer for even buffer sizes, and half minus 1 for odd buffer sizes, to be precise.

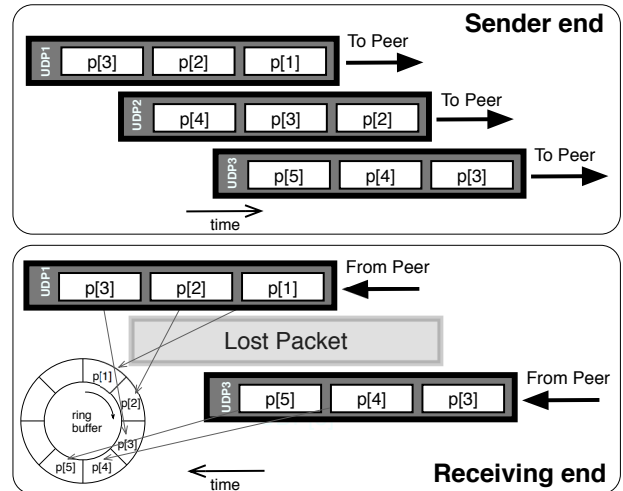


Figure 3. UDP redundancy,  $RedunFactor = 3$

---

### Algorithm 1 Packet redundancy receiving end

---

```

For every new UDP packet
  CurSeqNum ← Packet highest sequence num
  if (CurSeqNum - LastSeqNum) ≤ RedunFactor then
    NumNewPackets = CurSeqNum - LastSeqNum
  else
    NumNewPackets = RedunFactor
  end if
  for i = (NumNewPackets - 1) to 0 do
    Send P[CurSeqNum - i] to Ring Buffer
  end for
  LastSeqNum ← CurSeqNum

```

---

On the receiving end, the Algorithm 1 reads a UDP packet and determines if it has not already received the ex-

tra copies. New copies are sent to the ring buffer, and extra copies are discarded. Lost packets are recovered as illustrated in Figure 3.

## 2.5. Processing Plugins

JackTrip also has the ability to dynamically add plugins into the audio process callback (Fig. 1). One plugin implements loopback mode, i.e., audio received from a peer is immediately sent back. This allows a location to listen to its echo from a remote peer. The aural evaluation of network quality [6] or the synchronization of music through “feedback locking” [4] are two practical applications which use this approach.

Plugins can also be used for “Internet acoustics” or sonification through physical models [8], e.g., a network implementation of the Karplus-Strong algorithm for strings and drums synthesis.

## 3. CONCLUSIONS AND FUTURE WORK

“Broader” broad-band networks have the capacity to support high-quality audio. JackTrip serves to illustrate some of the software design decisions for achieving low-latency, bi-directional audio using these networks. Its particular uses have different requirements: collaborative music making is different from, for example, one-way remote studio recording where latency is not the issue but packet loss is. Depending on the application, JackTrip allows the user to tune its configuration, for example trading off some reliability (allowing for minor glitches) in favor of tighter latency.

At the design stage, the engineer must provide the methods that support the highest-quality musical performance. In particular, JackTrip deals with packet loss by providing a redundancy algorithm, and deals with clock drifts and late or un-recoverable packets by using a lower-level strategies in ring buffers that can, e.g., sound like a wavetable synthesizer, thus extending the musical sonority of the moment. Clock drift between remote WAN machines is still an unsolved issue and there are presumably new techniques to be tried in the future, like adaptive re-sampling, packet cross-fading, and others.

The current work of JackTrip is focused on the application layer, but new network projects like OpenFlow [1] and Dynamic Circuit Network [2] provide the opportunity to start experimenting with lower layers; it would be possible to dynamically specify network paths to minimize latency, or to obtain dedicated bandwidth for a more reliable Quality of Service (QoS).

## 4. ACKNOWLEDGMENTS

Fundamental contributors to design and coding have included Scott Wilson, Randal Leistikow and Daniel Walling.

At CCRMA, Fernando Lopez-Lezcano and Carr Wilkerson have provided continuous technical support and advice. Various individuals have also contributed in testing the software during several concerts, in particular Alain Renaud and Jonas Braasch.

## 5. REFERENCES

- [1] (2008) The OpenFlow Switch Consortium. [Online]. Available: <http://www.openflowswitch.org/>
- [2] (2009) Internet2 Dynamic Circuit Network. [Online]. Available: <http://www.internet2.edu/network/dc/>
- [3] (2009) JACK: Connecting a world of audio. [Online]. Available: <http://jackaudio.org/>
- [4] J.-P. Cáceres, R. Hamilton, D. Iyer, C. Chafe, and G. Wang, “To the edge with china: Explorations in network performance,” in *ARTECH 2008: Proceedings of the 4th International Conference on Digital Arts*, Porto, Portugal, 2008, pp. 61–66.
- [5] C. Chafe and M. Gurevich, “Network time delay and ensemble accuracy: Effects of latency, asymmetry,” in *Proceedings of the AES 117th Convention*, San Francisco, 2004.
- [6] C. Chafe and R. Leistikow, “Levels of temporal resolution in sonification of network performance,” in *Proceedings of the 2001 International Conference on Auditory Display*. Helsinki: ICAD, 2001.
- [7] C. Chafe, S. Wilson, R. Leistikow, D. Chisholm, and G. Scavone, “A simplified approach to high quality music and sound over IP,” in *Proceedings of the COST G-6 Conference on Digital Audio Effects (DAFX-00)*, Verona, Italy, Dec. 2000.
- [8] C. Chafe, S. Wilson, and D. Walling, “Physical model synthesis with application to internet acoustics,” in *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, Orlando, 2002.
- [9] D. E. Comer, *Internetworking with TCP/IP, Vol 1*, 5th ed. Prentice Hall, Jul. 2005.
- [10] S. Letz, Y. Orlarey, and D. Fober, “Jack audio server for multi-processor machines,” in *Proceedings of International Computer Music Conference*, ICMA, Ed., Barcelona, 2005, pp. 1–4.
- [11] A. Xu and J. R. Cooperstock, “Real-time streaming of multichannel audio data over Internet,” in *Proceedings of the 108th Convention of the Audio Engineering Society*, Paris, 2000, pp. 627–641.