# Flexilink: A unified low latency network architecture for multichannel live audio

Yonghao Wang[1], John Grant[2], and Jeremy Foss[3]

[1] Birmingham City University, Birmingham, B4 7XG, UK
yonghao.wang@bcu.ac.uk

[2] Nine Tiles Networks Ltd, Cambridge, CB25 9HT, UK
j@ninetiles.com

[3] Birmingham City University, Birmingham, B4 7XG, UK
jeremy.foss@bcu.ac.uk

## ABSTRACT

The networking of live audio for professional applications typically uses layer 2 based solutions such as AES50 [1] and MADI utilising fixed time slots similar to Time Division Multiplexing (TDM).  However, these solutions are not effective for best effort traffic where data traffic utilises available bandwidth and is consequently subject to variations in QoS. There are audio networking methods such as AES47 which is based on asynchronous transfer mode (ATM), but ATM equipment is rarely available. Audio can also be sent over Internet Protocol (IP), but the size of the packet headers and the difficulty of keeping latency within acceptable limits make it unsuitable for many applications. In this paper, we propose a new unified low latency network architecture that supports both time deterministic and best effort traffic towards full bandwidth utilisation with high performance routing/switching. For live audio, this network architecture allows low latency as well as the flexibility to support multiplexing multiple channels with different sampling rates and word lengths.

## 1.  BACKGROUND

In a digital system supporting interactive media applications such as live audio, the end-to-end latency is preferably measurable and manageable. For some applications such as in-ear monitoring, acceptable latency varies and can be as short as less than 2ms [2].

In an audio processing system such as in Figure 1, the end-to-end latency arises from different sources [3]: conversion from analogue to digital and back to analogue (ADC/DAC) [4]; networking and routing; the digital console; and the computer system with software plugins (DAW) [5] etc. Buffering is the major cause of latency in IP networks [6], especially in the Internet, but also including intranets. For professional live audio

applications, where low latency is required, closed networks with the audio specific layer 1 and layer 2 technology are commonly adopted.
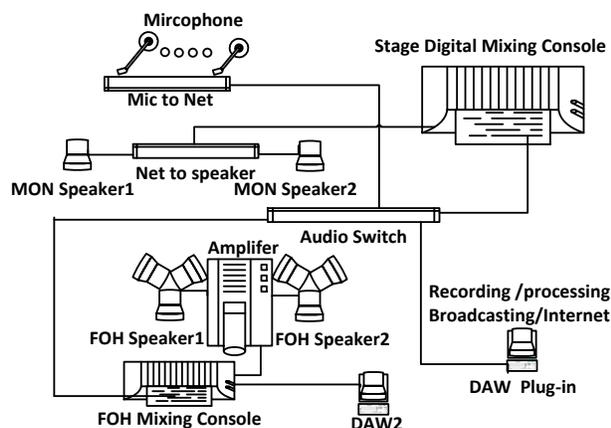


Figure 1 Audio Processing System

In high resolution and low latency audio applications, many audio specific networking technologies modify the existing layer-2 or layer-3 protocols to utilise the current lower layer network infrastructure such as Ethernet. However, it is difficult to achieve a networking architecture, which supports both time critical audio data and also best effort data (such as file transfer and emails). Converged networks based on IP, scalable from LAN to WAN are required to support the vast (and growing) interactive audio/video media traffic on the Internet. However connectionless packet architectures are inevitably a problem for deterministic data especially when low latency is required. Current QoS, traffic engineering and over-provisioning solutions cannot solve the entire problem: they complicate the system, and increase power requirements and cost. The professional audio networking industry is reluctant to use current Internet solution for time critical applications. Instead, they normally adopt solutions based on specific protocols designed and modified from physical layer up to layer 3 such as AES50, EtherSound, CobraNet etc.  For low latency live audio, TDM based protocols such as AES50 can provide very good performance. The proprietary AES50 router can provide latency as low as a few samples with a fixed number of channels reserved for packet data.

However, it appears that there is no unified network solution to provide flexible and bandwidth efficient support for both low latency deterministic traffic and best effort traffic. There is also an issue with multi-

channel digital audio streams with a range of sampling frequencies and variable bit lengths, and the need for flexible routing and channel assignments. Current multiplexing methods are insufficient to support them without sampling rate conversion and data format rectification.

## 2.   INTRODUCTION

The proposed architecture is to effectively support both best effort data and time deterministic data (audio sample packets). It should also interwork with existing network infrastructure and protocols at maximum compatibility. Since the current physical network layer (such as full duplex Ethernet) can be viewed as time deterministic bit pipe  there is no reason why a time deterministic logical control layer cannot be implemented. This allows a guaranteed Quality of Service (QoS) and expected Quality of Experience (QoE) for the higher layer protocols. The exact time delay for the transmission of time-critical data can therefore be estimated.

Based on earlier work [7], we proposed a novel unified network architecture that combines the advantages of TDM and best effort networks. The proposed layer 2 protocols, "Flexilink", have been developed along with a prototyped network processor architecture and interface cards. Compatibility with existing Ethernet infrastructure is maintained. Flexilink can operate at full-duplex mode, where non-deterministic CSMA-CD can be avoided.

## 3.   THE ARCHITECTURE DESIGN OF FLEXILINK

### 3.1.   The rationale of the design

User generated data can be categorised as (i) data to be transmitted as time-deterministic with constant intervals and predicable delay, i.e. real-time data; (ii) data to be transmitted at earliest opportunity but without the constraints of real time, i.e. best effort data. The network also conveys network management data.

This gives us the three data categories as follows:

*   Synchronous flow (SF) for audio/video and other time deterministic data.

- Asynchronous flow (AF) for best effort data.

- Control Message (CM) for session control and link management.

The theoretical requirements of a single SF can be determined; for example - transmitting a 44.1kHz sampled CD with 16 bit samples, without any headers and error checking mechanism will require 1.4112Mb/s. Compressed formats will also have a nominal bit rate allocated (with the associated compromised quality). So, for a link with sufficient bandwidth, we are able preallocate spaces or slots for the SF data packets. AF data can be transmitted in the gaps between SF data packets. A simple theoretical link model is shown in Figure **2**, where SF data packets are transmitted at constant time $t_0$. Since SF packets are of variable length, gaps to be filled are also variable.
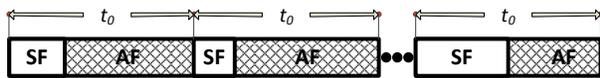


Figure 2 Ideal Link for the Traffic

To ensure SF data are transmitted in a time deterministic manner, resources (bandwidth requirements) need to be reserved when the link is established. Individual SF data packets are identified by the position of the SF packets within the stream (similar to time slots in a TDM frame).

Control messages (CM) with associated protocols are used for establishing links and negotiating resource reservations.

Flexilink supports variable length SF packets where the varying length gaps are filled by any AF traffic awaiting transmission. To facilitate this a small header is added to an SF packet. The header is simplified to contain only the length of the SF packet plus basic error checking bits. Therefore there is no need for AF data to be encapsulated with a new header when it is fragmented by the SF flow. This also simplifies the hardware logic required to forward both SF and AF traffic effectively.

This operation could be considered as a continuous AF stream frequently interrupted by the frequent real-time SF data, since the main parameters are all known: the speed of network link, the data rate of SF, and length of the SF data packets. There is no additional reassembly required to reconstruct the segmented AF traffic. In addition, this design can achieve the maximum

utilisation of link bandwidth with all the gaps (unused capacity) filled by the available best effort traffic. [8] and [9] proposed a similar system, but [8] has fixed TDM channels so that the capacity allocated but not carrying data is not utilised. [9] has two types of traffic, but the low priority traffic is fragmented with an additional header carrying type and destination information being required for each fragment.

## 3.2. Architecture Design

The network node supporting the proposed Flexilink would have a common network architecture below in figure 3, with two major functional blocks: the "control unit" for setting up and tearing down call flows, allocating the resources, and route finding algorithms; and the data "forwarding logic" for fast forwarding and switching the data for both SF and AF data.
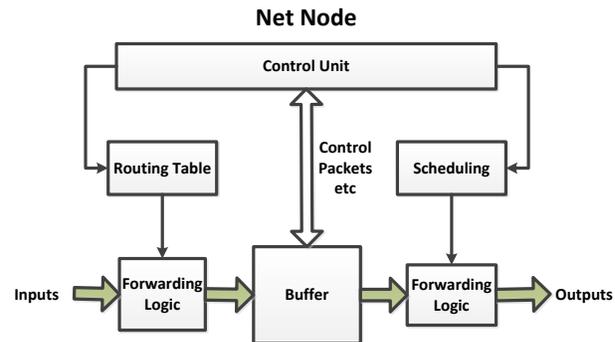


Figure 3 Architecture of Network Node

First of all, there is a process for setting up and tearing down the link between two end points with an intelligent time slot map allocation algorithm, which is based on the available network link resources and the requirements of deterministic traffic (SF). This setup can be managed by control messages (CM).

The CM can be implemented as standard IEC 62379-5-2 (Common Control Interface for networked audio and video products) messages [10]. Essentially they are considered as normal AF traffic for the purposes of the link. However, whereas normal AF traffic is routed to the output to be transmitted over the link, the CMs are directed to the controller. CMs have priority over AFs on each link.

For audio traffic, the packages in a SF can be as small as audio samples, for example 48000 packets per second with a payload of 4 bytes.

### 3.2.1. Design Header of SF

The simple header added to SFs contain only the length information. To minimise the header cost, the length of the header is also variable as shown in Figure 4 (a):

- 1-byte header to support 0 up to 15 bytes SF packet length.

- 2-byte header to support 16 up to 255 bytes SF packet length.

- 3-byte header to support 256 up to 4095 bytes SF packet length.

A 1-byte header consists of 4 bits to encode the length information; 3 bits for CRC; and 1 bit flag to indicate if there are further header bytes as shown in Figure 4 (b).
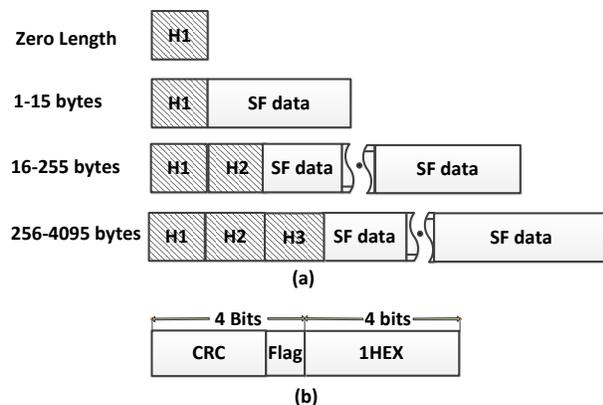
Figure 4 Header of SF packet

### 3.2.2. Interface architecture to support Flexilink

To maximise compatibility, Flexilink should be able to use the existing physical network interface. However to support the proposed Flexilink protocol, a new media access control (MAC) layer architecture needs to be considered, which allows AF and SF to be treated differently. Figure **5** shows the simplified Flexilink MAC layer in which AF and SF have separate buffers and copy logic allocated for them.
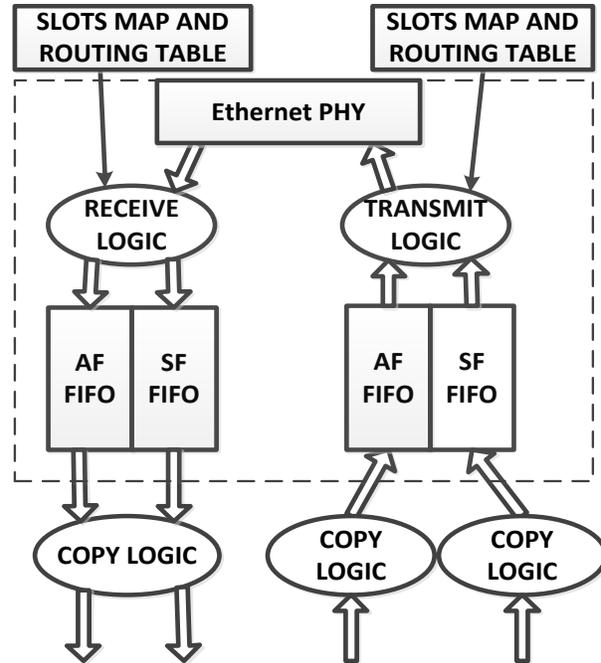
Figure 5 Simplified Flexilink MAC layer

### 3.2.3. Layered Traffic Model

The theoretical link traffic slot allocation is shown in Figure 2. The SF and AF access may be implemented over an existing point-to-point link mechanism in order to utilise current network infrastructure. Figure 6 shows a practical implementation of the layered traffic model.

Figure 6 Layered Flexilink Traffic Model

The Frame layer can be a standard fixed size Ethernet frame, so the position of an SF packet in relation to the start symbol of frame can be used as a reference for identification of the SF data packets.

### 3.2.4. Supporting flexible multichannel audio streams with different sampling frequencies

A current problem in networking audio traffic is the lack of flexibility to support arbitrary numbers of audio channels with different sampling frequencies and of compressed or uncompressed data format. In the proposal given here, multiple channels of different

sampling frequencies can be readily supported as long as the capacity of the link above the frame layer is greater than the total bandwidth requirements of number of the SFs.

For Ethernet physical media, the Ethernet Jumbo Frame format can be used to maximise the capacity of available bandwidth and minimise the cost of inter-frame gaps. The time slot map is allocated within the payloads of the frames. The time slot map allocation algorithm will ensure the SF streams are distributed evenly. A phase control algorithm corrects the delivery of samples at end points to ensure the samples are rendered at certain jitter requirements. The Precision Time Protocol can be adopted for accurate timing and provision of clocking.

## 4. THE ETHERNET IMPLEMENTATIONS OF FLEXILINK

The Flexilink design can be implemented at different types of physical layer as long as the bit transfer rate can be fixed and the bit pipe is guaranteed.

This section will discuss how Flexilink is implemented using 1 Gigabit Ethernet infrastructure, since 90% internet traffic is originated from Ethernet or WiFi, and Gigabit Ethernet is commonly used by other proprietary audio network technology.

### 4.1. Ethernet Frame structure for Flexilink

To maximise bandwidth efficiency, Flexilink adopts the Jumbo Frame format in which the payload size is greater than the standard 1500bytes common Ethernet frame. The maximum size of a Jumbo frame cannot exceed 11455 bytes in order to allow CRC algorithm to effectively working [11]. Another limitation is that some Ethernet PHY chips do not support frames larger than 9000 bytes.

To be able to accurately access the allocation slots and the positions of SF packets, it is preferred to have a fixed allocation period at the frame layer. In this case the allocation period of 124.96µs is chosen to be slightly faster than 125µs since the 125µs is the frame cycle used in many other designs such as Synchronous Digital Hierarchy (SDH), Firewire and full speed USB.

The network devices can use standard AES51 [12] negotiation packets (which are standard Ethernet MAC

packets) to setup the links. Once both ends are in Flexilink mode, a "Reduced Jumbo Frame" (RJF) format is utilised to maximise the payload. The RJF eliminates some unnecessary parts of the standard Ethernet frame (for example, the source and destination addresses, because all frames are sent from the node at one end of the link to the node at the other) to give more payload space to the allocation periods.

 It is described as below:

- 2 bytes preamble + Start Frame Delimiter (SFD).

- 5 bytes AES51 packet type and timing information.

- 7785 bytes payload data.

- 4 bytes FCS.

- 14 bytes inter frame gap (IFG).

In total the RJF frame size including IFG is 7810 bytes long. Two successive RJF combine together to make a 124.96µs allocation period (AP) at full duplex 1 Gigabit Ethernet link. This design is to guarantee the 8000 allocation periods/second can be transmitted over 1 Gigabit Ethernet link. Each allocation period has 15570 bytes payload space to transmit SF and AF traffic. The theoretical bandwidth utilisation can up to 99.6%.

The FCS is only used to check that the link is working reliably. Routing of AFs and SFs does not wait until the FCS has been received, and for many media formats it is better to deliver data with a few bit errors than to discard whole frames.

### 4.2. Methods for support low latency multiple sampling rate audio streams

The following example demonstrates how Flexilink supports multichannel audio streams with different sampling frequencies.

Assuming that we have audio streams as (i) Flow 1: 48kHz mono 24bits; (ii) Flow 2: 44.1kHz stereo 16bits; and (iii) Flow 3: 96kHz mono 24 bits; we need to transmit (or multiplex) them using Flexilink.

The audio data of all three flows plus other possible metadata is less than 15 bytes; therefore 1 byte SF header is needed for each packet.

The control unit of the sender node will reserve the network resource and allocate transmitting slots for each flow as SF packets; i.e. the slot allocation map will be established and agreed by both ends of the link by CM messages. Within each packet, in addition to the audio data, there is 1 byte containing synchronisation information as specified in 7.3.2 of IEC 62379-5-2 [10]. (An additional 1 byte per channel of overhead may also be added to carry channel status CRC, etc.)

- For flow 1: The number of allocation slots within allocation period should be $\geq$ 48/8. Therefore 6 slots are needed from one allocation period with each slot being 3+1+1 = 5 bytes long.
- For flow 2: The number of allocation slots $\geq$44.1/8 = 5.5. Hence 6 slots as well, with each slot being 4+1+1 = 6 bytes long.
- For flow 3: The number of allocation slots $\geq$ 96/8. Hence also 12 slots, with each slot being 3+1+1 = 5 bytes long.

Note that for flow 2 although the number of slots we allocated is more than actually needed for delivery of audio samples at 44.1Khz, this is because the number of slots is rounded up to a whole number. Approximately every twelfth slot there will be empty SF data packet for flow 2. However this is not a problem since the header will indicate an empty packet, therefore the space can be used for AF data.

Having AES51 packet type and timing information in the Flexilink over Ethernet implementation, the sender and receiver can be synchronised easily. When accurate time information needs to be distributed, the (Precision Time Protocol) PTP can also be used for synchronisation between sender and receiver to avoid potential drift of clock and jitter over a period of transmission.

Typical latency for synchronous flows would be 3 to 6 microseconds per hop, plus the "speed of light" delay in the transmission line.

## 5. HARDWARE IMPLEMENTATION

The Flexilink architecture needs be implemented on various network devices such as single port or multiport interface cards, and network routers or switches. The key building blocks of these devices are the data forwarding unit and the control unit.

The data forwarding unit needs to fast forward uses cut-through forwarding for the SF data packets as well as storing and forwarding the AF data, which can be implemented in hardware logic by FPGA. The control unit needs to flexibly allocate the resource for multichannel SF data requests, such as positioning the SF data packet in the continuous bit stream, and allocating the SF data packets evenly to allow phase correction algorithms work efficiently. The control unit needs to be an effective general purpose CPU with accelerated networking processing capabilities.

At Nine Tiles Networks Ltd, a prototype control unit has been built on an XC6S LX9 FPGA, dRAM and flash memory.



Figure 7 Picture of Prototype board

The processor (which is implemented in the FPGA fabric) and the high level language in which it is programmed are optimised for processing protocol messages.

## 6. INTERCONNECTION WITH EXISTING NETWORK

In a mixed network environment multiple low latency audio flows of different sampling frequencies, bit depths and numbers of channels, can be well supported by Flexilink as described in section 4.2,

Normal data transfer, typically IP-based traffic (e.g. file transfer), is mapped into the AF traffic and so is safely transmitted over a Flexilink network.

The architecture design maintains separation between the AF and SF data and thus ensures that there is no interference between two types of data whilst fully utilising the bandwidth not used by SF data.

For IP based audio data, Flexilink can map audio IP packets to SF data packet according to the identified

service priority. So Flexilink should not negate the original QoS.

In the case where a Flexilink interface is peering with a normal Ethernet interface which does not support the Flexilink mode, Flexilink can (i) switch to standard Ethernet mode or (ii) negotiate an Ethernet AVB mode to prioritize delivery of the audio traffic.

## 7. FUTURE WORK AND IMPACT

Apart from the audio application discussed here, Flexilink can be a good candidate layer 2 technology to support guaranteed QoS for upper layer applications, such as synchronised audio/video delivery and low latency interactive media distribution networks.

The Flexilink network architecture also provides an ideal solution for a truly QoS guaranteed end-to-end Integrated Service, although more development is required on the scalability and interoperability with current rapidly evolving networks.

## 8. ACKNOWLEDGEMENTS

## 9. REFERENCES

[1] Audio Engineering Society, 'AES50-2011 AES standard for digital audio engineering - High-resolution multi-channel audio interconnection (HRMAI)', *Audio Engineering Society, Inc.*, 2011.

[2] M. Lester and J. Boley, 'The Effects of Latency on Live Sound Monitoring', in *Audio Engineering Society Convention 123*, 2007.

[3] N. Bouillot, E. Cohen, J. R. Cooperstock, A. Floros, N. Fonseca, R. Foss, M. Goodman, J. Grant, K. Gross, S. Harris, and others, 'AES White Paper: Best Practices in Network Audio', *J. Audio Eng. Soc*, vol. 57, no. 9, p. 729, 2009.

[4] Y. Wang, 'Engineering Brief: "Latency Measurements of Audio Sigma Delta Analogue to Digital and Digital to Analogue Converts "', presented at the 131st AES Convention, New York, NY, USA, 2011.

[5] Y. Wang, R. Stables, and J. Reiss, 'Audio Latency Measurement for Desktop Operating Systems with Onboard Soundcards', in *Audio Engineering Society Convention 128*, 2010.

[6] J. P. Cáceres and C. Chafe, 'JackTrip: Under The Hood Of An Engine For Network Audio', *Journal of New Music Research*, vol. 39, no. 3, pp. 183–187, 2010.

[7] J. S. Grant, 'Method And Apparatus For Transceiving Data', 23-Jul-2010. [Online]. Available: http://patentscope.wipo.int/search/en/WO20100820 42. [Accessed: 20-Jul-2012].

[8] Dennis Gordon Froggatt, 'Data Transmission System', 22-Jan-1986.

[9] P. Strong, T. Wild, and G. Dean, 'Latency Reduction By Adaptive Packet Fragmentation', 07-Mar-2008.

[10] IEC Project Team 62379, 'IEC 62379 Common Control Interface For Networked Digital Audio And Video Products - Part 5-2: Transmission Over Networks - Signalling'. Draft, 2012.

[11] Alteon Networks, 'Extended Frame Sizes for Next Generation Ethernets'. 1998.

[12] Audio Engineering Society, 'AES51-2006 (r2011) AES Standard for Digital Audio – Digital Input-Output Interfacing – Transmission of ATM Cells Over Ethernet Physical Layer', *Audio Engineering Society, Inc.*, 2006.