



**IMT Atlantique**

Bretagne-Pays de la Loire

École Mines-Télécom

# *INTRODUCTION À PSI*

Cédric Dumas & Aurélien Milliat

IHM 2024 : La collaboration en présence mixte dans les espaces interactifs

# SOMMAIRE

## 1. PSI

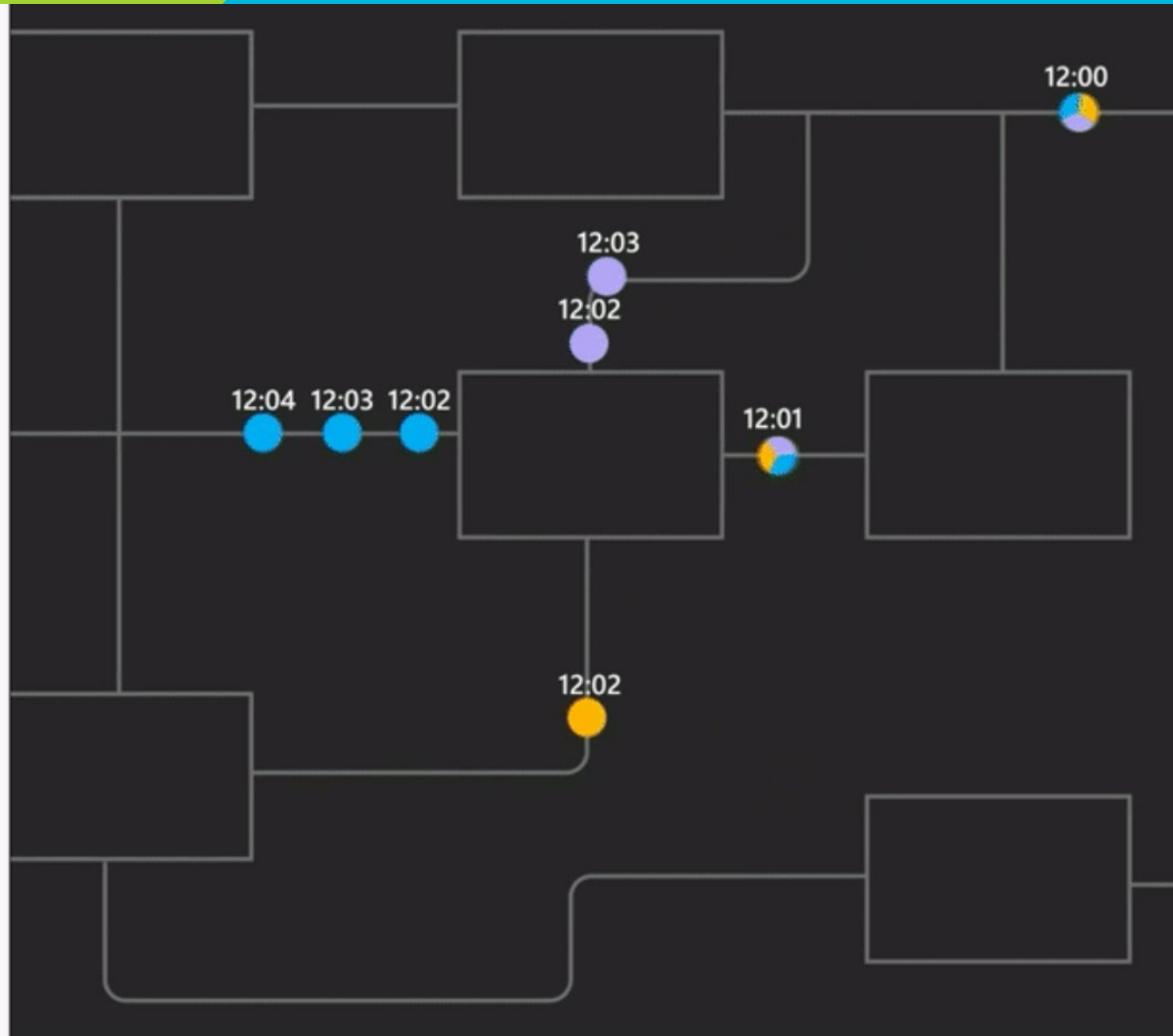
- i. PRÉSENTATION
- ii. CAS D'USAGES
- iii. AUTRES SOLUTIONS
- iv. BILAN
- v. PSI STUDIO
- vi. PROGRAMMATION
- vii. EXEMPLE

## 2. UN PEU DE PRATIQUE

- i. HELLO WORLD
- ii. WEBCAMWITHAUDIO & STORE
- iii. KINECT AZURE
- iv. OFFLINE PROCESSING
- v. \PSI, WEBRTC ET UNITY



# INFRASTRUCTURE FOR SYNCHRONIZATION



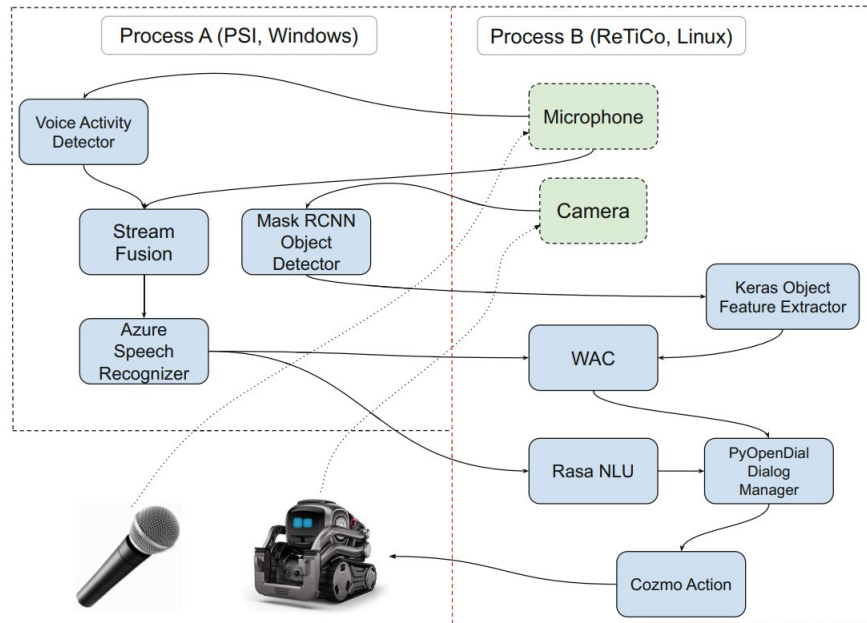
**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

<https://camo.githubusercontent.com/1f2c3afc838e595a3d8738c25877b9cb108fd5576eef23e543e4a2985ab1ae11/68747470733a2f2f777772e6d6963726f736f66742e636f6d2f656e2d75732f72657365617263682f75706c6f6164732f70726f642f323031382f30312f5073692d476966322d3936302d436f727265637465642e676966>

# Platform for Situating Intelligence Overview

- ▶ Plateforme de synchronisation de flux de données.
- ▶ Principe de composants / pipeline en open source (MIT License).
- ▶ Possibilité d'une utilisation réseau / remote (Hololens).
- ▶ Capacité d'enregistrement des flux.
- ▶ Possibilité de temps réel (ou de simuler via les enregistrements).
- ▶ Orienté Machine Learning
- ▶ Application modulaire de visualisation: Psi Studio.

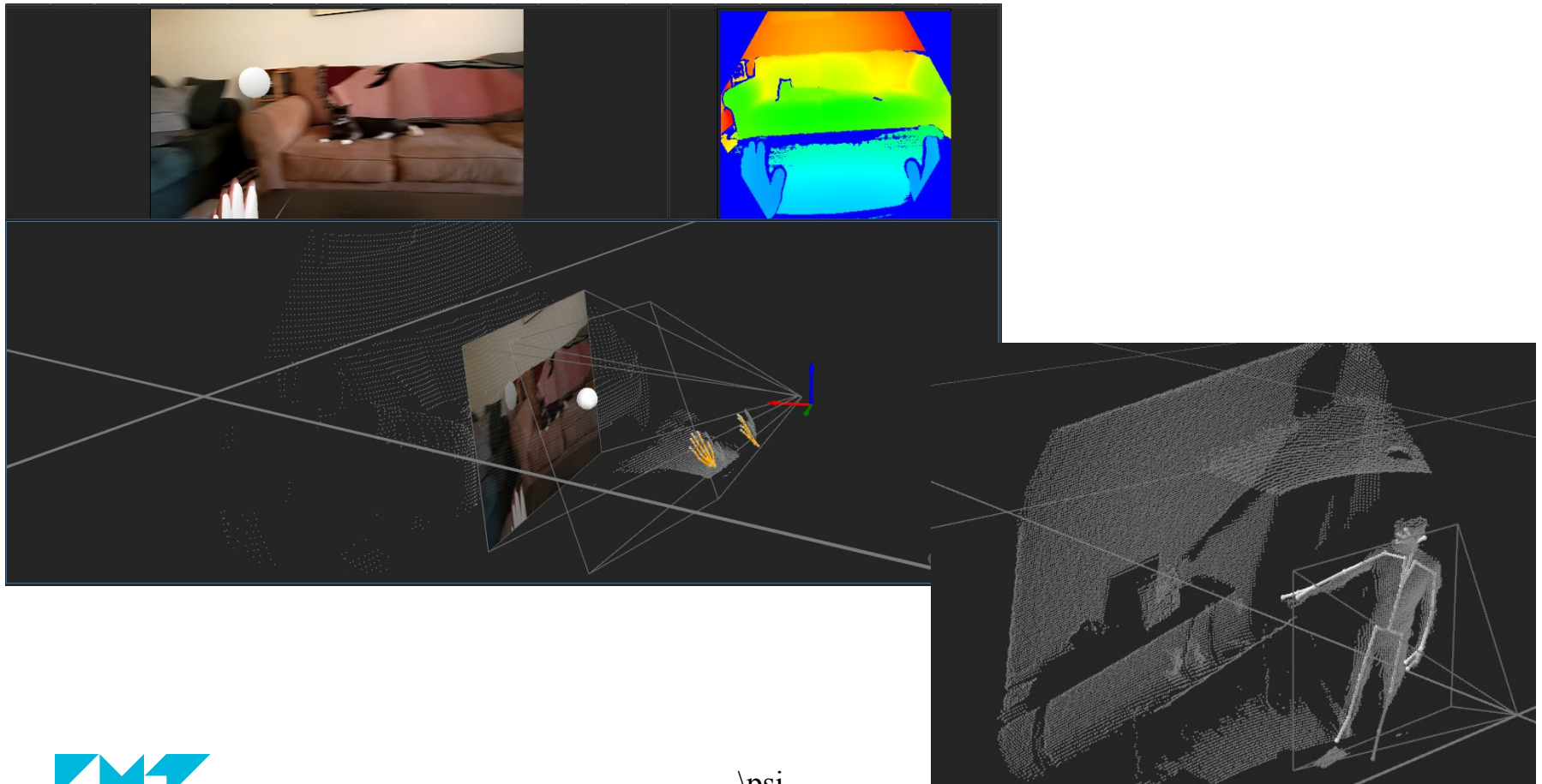
## ► Human Robot Interaction (HRI) :



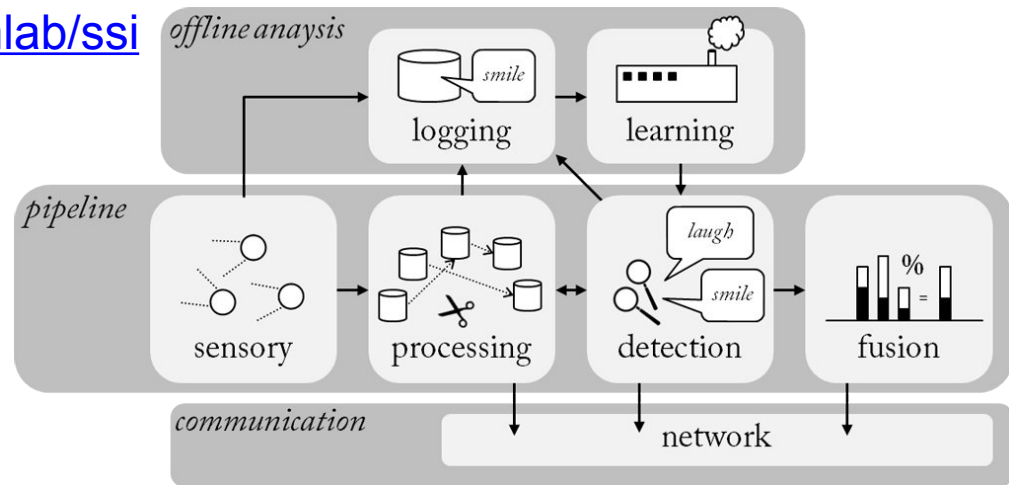
<https://link.springer.com/content/pdf/10.1007/978-3-031-06018-2.pdf>

► **exemple : HARMONI** (Composing HARMONI: An Open-source Tool for Human and Robot Modular Open Interaction)

## ► Mixed Reality & Collaborative Interaction:



- ▶ Social Signal Interpretation (SSI) Framework:
  - ▶ Human Centered Artificial Intelligence Lab of the Augsburg University
  - ▶ Beaucoup de composants implémentés : capteurs et libraires de process.
  - ▶ Principalement en C++
  - ▶ Latest Release v1.0.4 17 Apr 2018
  - ▶ <https://github.com/hcmlab/ssi>



\psi



## ▶ OpenSense:

- ▶ Intelligent Human Perception Lab, USC-ICT
- ▶ Surcouche graphique à PSI.
- ▶ Composants pour OpenFace, OpenPose, OpenSmile & Whisper.
- ▶ Possibilité de constituer un pipeline de traitement en 'no code'.
- ▶ Pas mal de bugs.
- ▶ IHM pas très 'user friendly'.
- ▶ <https://github.com/intelligent-human-perception-laboratory/OpenSense/>

## ▶ MediaPipe:

- ▶ Développé par Google
- ▶ Orienté traitement d'image -> segmentation d'image.
- ▶ Mais framework générique.
- ▶ Disponible sur plusieurs plateformes
- ▶ <https://github.com/google/mediapipe>

	Android	iOS	C++	Python	JS	Coral
Face Detection	✓	✓	✓	✓	✓	✓
Face Mesh	✓	✓	✓	✓	✓	
Iris	✓	✓	✓			
Hands	✓	✓	✓	✓	✓	
Pose	✓	✓	✓	✓	✓	
Holistic	✓	✓	✓	✓	✓	
Selfie Segmentation	✓	✓	✓	✓	✓	
Hair Segmentation	✓		✓			
Object Detection	✓	✓	✓			✓
Box Tracking	✓	✓	✓			
Instant Motion Tracking	✓					
Objectron	✓		✓	✓	✓	
KNIFT	✓					
AutoFlip			✓			
MediaSequence			✓			
YouTube 8M			✓			

## PROS:

- ▶ C#
- ▶ Relativement simple
- ▶ Store/Remote
- ▶ Microsoft
- ▶ Lien avec le cloud Azure mais ouvert
- ▶ Composants existants
- ▶ Compatible Unix
- ▶ Linq

## CONS:

- ▶ Performance/ressource (bench?)
- ▶ In Dev (certaines limitations)
- ▶ Pas de plugin Unity/Unreal

► Tarif des services Azure: <https://azure.microsoft.com/fr-fr/pricing/details/cognitive-services/> (11-2023)

Produit	Fonctionnalités	Tarif
Compréhension du langage (LUIS) Standard jusqu'à 50 requêtes par seconde	Demandes de texte	<b>1,418 €</b> toutes les 1 000 transactions* de prédiction
	Demandes vocales	<b>5,199 €</b> toutes les 1 000 transactions
Langage Standard jusqu'à 1 000 requêtes par seconde et 1 000 requêtes par minute	Analyse des sentiments (et exploration des opinions)	0millions-0.5millions d'enregistrements de texte - <b>0,9453 €</b> par 1 000 enregistrements texte
	Extraction des expressions clés	0.5millions-2.5millions d'enregistrements de texte -
	Détection de la langue	<b>0,7090 €</b> par 1 000 enregistrements texte
	Reconnaissance d'entité nommée, y compris les informations d'identification personnelle (non disponible en conteneur)	2.5millions-10.0millions d'enregistrements de texte - <b>0,2836 €</b> par 1 000 enregistrements texte Enregistrements de texte 10.0M+ - <b>0,2364 €</b> par 1 000 enregistrements texte
	Analyse de texte pour le secteur de la santé	0millions-0.005millions d'enregistrements de texte - Inlus
		0.005millions-0.5millions d'enregistrements de texte - <b>23,6318 €</b> par 1 000 enregistrements texte
		0.5millions-2.5millions d'enregistrements de texte - <b>14,1791 €</b> par 1 000 enregistrements texte Enregistrements de texte 2.5M+ - <b>9,4527 €</b> par 1 000 enregistrements texte
	Réponses aux questions*	Gratuit
Traducteur S1	Traduction de texte Personnalisation Détection de la langue Dictionnaire bilingue Translittération	<b>9,453 €</b> par million de caractères (Paiement à l'utilisation)

# PSI - Bilan

Produit	Fonctionnalités	Tarif
Vision par ordinateur S1 30 transactions par seconde pour les opérations de lecture 15 transactions par seconde pour les opérations non-lecture	Balise Visage GetThumbnail Couleur Type d'image GetAreaOfInterest Détection de personnes* Cultures intelligentes	0-1 million de transactions - <b>0,946 €</b> toutes les 1 000 transactions 1-10 million de transactions - <b>0,615 €</b> toutes les 1 000 transactions 10-100 million de transactions - <b>0,568 €</b> toutes les 1 000 transactions +100 million de transactions - <b>0,379 €</b> toutes les 1 000 transactions
	OCR Adulte Célébrité Repère Détecter, Objets Marque	0-1 million de transactions - <b>0,946 €</b> toutes les 1 000 transactions 1-10 million de transactions - <b>0,615 €</b> toutes les 1 000 transactions 10-100 million de transactions - <b>0,568 €</b> toutes les 1 000 transactions +100 million de transactions - <b>0,379 €</b> toutes les 1 000 transactions
	Décrire* Lecture Sous-titre*	0-1 million de transactions - <b>1,418 €</b> toutes les 1 000 transactions +1 million de transactions - <b>0,568 €</b> toutes les 1 000 transactions
Content Moderator S0 jusqu'à 10 requêtes par seconde	Modérées, Révision	0-1 million de transactions - <b>0,946 €</b> toutes les 1 000 transactions 1-5 million de transactions - <b>0,709 €</b> toutes les 1 000 transactions 5-10 million de transactions - <b>0,568 €</b> toutes les 1 000 transactions +10 million de transactions - <b>0,379 €</b> toutes les 1 000 transactions
API Face Standard jusqu'à 10 requêtes par seconde	Détection de visage Vérification de visage Identification de visage Regroupement de visages Recherche de visages semblables	0-1 million de transactions - <b>0,946 €</b> toutes les 1 000 transactions 1-5 million de transactions - <b>0,757 €</b> toutes les 1 000 transactions 5-100 million de transactions - <b>0,568 €</b> toutes les 1 000 transactions +100 million de transactions - <b>0,379 €</b> toutes les 1 000 transactions
	Face Storage	<b>0,010 €</b> par 1 000 visages par mois

The screenshot displays the PSI Studio interface, titled "Platform for Situated Intelligence Studio - Untitled Dataset". The main visualization area shows a 3D scene with a person's skeleton and a camera view of a room containing a television, a plant, and a table. The interface includes a "Datasets" panel on the left, a "Layout" panel, and a "Properties" panel on the right. The "Properties" panel is set to "Audio.VAD" with a color of "DodgerBlue".

**Layout Panel:**

- Timeline Panel
  - Pointing.IsPointing
- Timeline Panel
  - Sources.Audio
  - Audio.VAD
- Timeline Panel
  - Audio.Features.LogEnergy
- Timeline Panel
  - Audio.Speech.Results
- Timeline Panel
  - Pointing.VisionResults (Latency)
- Instant Panel
  - 2D Panel
    - Pointing.ElbowAngle.Histogram
  - 2D Panel
    - Pointing.HandTipLifted.Histogram

**Properties Panel (Audio.VAD):**

- Color: DodgerBlue
- Interpolation Style: Step
- Legend Format: 0
- Legend Value: 0
- Line Width: 2
- Marker Color: LightGray
- Marker Size: 4
- Marker Style: None
- Name: Audio.VAD
- Range Color: DodgerBlue
- Range Line Width: 1
- Stream Adapter: BoolToDoubleAdapter
- Summarizer: DoubleRangeSummarizer
- Visible:
- Y Axis Compute Mode: Manual
- Y Max: 1
- Y Min: 0

**Visualizations:**

- Timeline: Shows audio waveforms and speech results with labels "what is that", "how about that", and "what's that one".
- Elbow Angle Histogram: Shows a distribution of elbow angles in degrees, with a peak around 137.25.
- Hand Tip Lifted Histogram: Shows a distribution of hand tip lifted values in Delta, with a peak around -0.4825.

- ▶ Visualisation des données (temps réel ou enregistrement).
- ▶ Création de dataset à partir des stores.
- ▶ Annotation !
- ▶ Analyse de données via API :
  - ▶ Exécution d'algorithme spécifique
  - ▶ Processus d'extraction pour Machine Learning
  - ▶ Statistiques sur les données
- ▶ Possibilité de développer des modules de visualisation pour l'application.

```
using (var p = Pipeline.Create())
{
    // Create the microphone, acoustic feature extractor, and connect them
    //var microphone = new AudioCapture(p, WaveFormat.Create16kHz1Channel16BitPcm());
    var inputStore = PsiStore.Open(p, "SimpleVAD", Path.Combine(Directory.GetCurrentDirectory(), "Stores", "SimpleVAD.0000"));
    var microphone = inputStore.OpenStream<AudioBuffer>("Audio");
    var acousticFeaturesExtractor = new AcousticFeaturesExtractor(p);
    microphone.PipeTo(acousticFeaturesExtractor);

    // Display the log energy
    acousticFeaturesExtractor.LogEnergy
        .Sample(TimeSpan.FromSeconds(0.2))
        .Do(logEnergy => Console.WriteLine($"LogEnergy = {logEnergy}"));

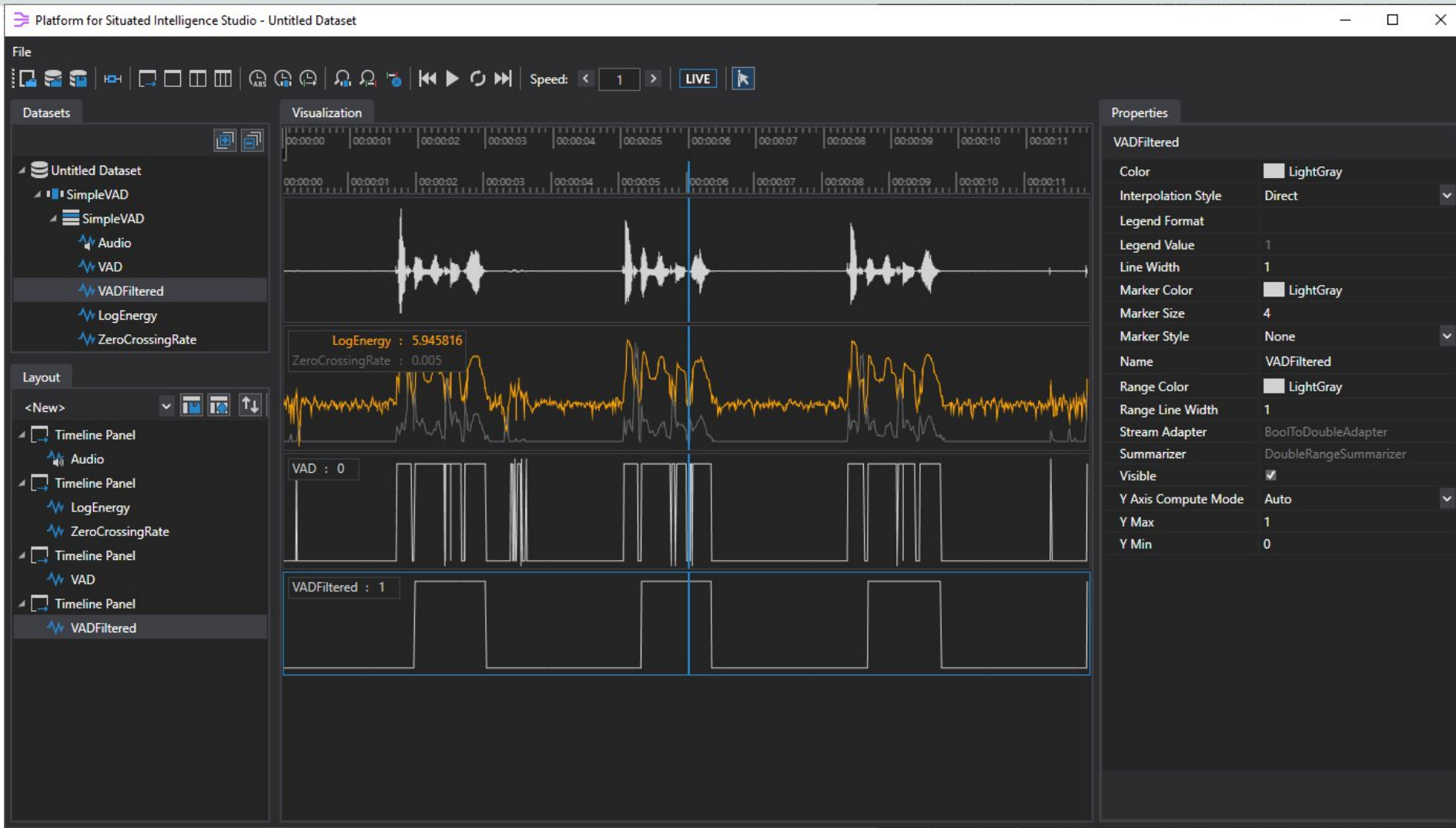
    // Create a voice-activity stream by thresholding the log energy
    var vad = acousticFeaturesExtractor.LogEnergy
        .Select(l => l > 7);

    // Create filtered signal by aggregating over historical buffers
    var vadWithHistory = acousticFeaturesExtractor.LogEnergy
        .Window(RelativeTimeInterval.Future(TimeSpan.FromMilliseconds(300)))
        .Aggregate(false, (previous, buffer) => (!previous && buffer.All(v => v > 7)) || (previous && !buffer.All(v => v < 7)));

    // Write the microphone output, VAD streams, and some acoustic features to the store
    var store = PsiStore.Create(p, "SimpleVAD", Path.Combine(Directory.GetCurrentDirectory(), "Stores"));
    microphone.Write("Audio", store);
    vad.Write("VAD", store);
    vadWithHistory.Write("VADFiltered", store);
    acousticFeaturesExtractor.LogEnergy.Write("LogEnergy", store);
    acousticFeaturesExtractor.ZeroCrossingRate.Write("ZeroCrossingRate", store);

    p.RunAsync();
    Console.ReadKey();
}
```





```
namespace ComponentNameSpace
{
    using System;
    using Microsoft.Psi;
    using Microsoft.Psi.Components;

    public class Component : Subpipeline
    {
        // Connector for the reciever
        protected Connector<TYPE_IN> InConnector;

        // Reciever that encapsulates the input.
        public Receiver<TYPE_IN> In => InConnector.In;

        // Emitter taht encapsulates the output.
        public Emitter<TYPE_OUT> Out { get; private set; }

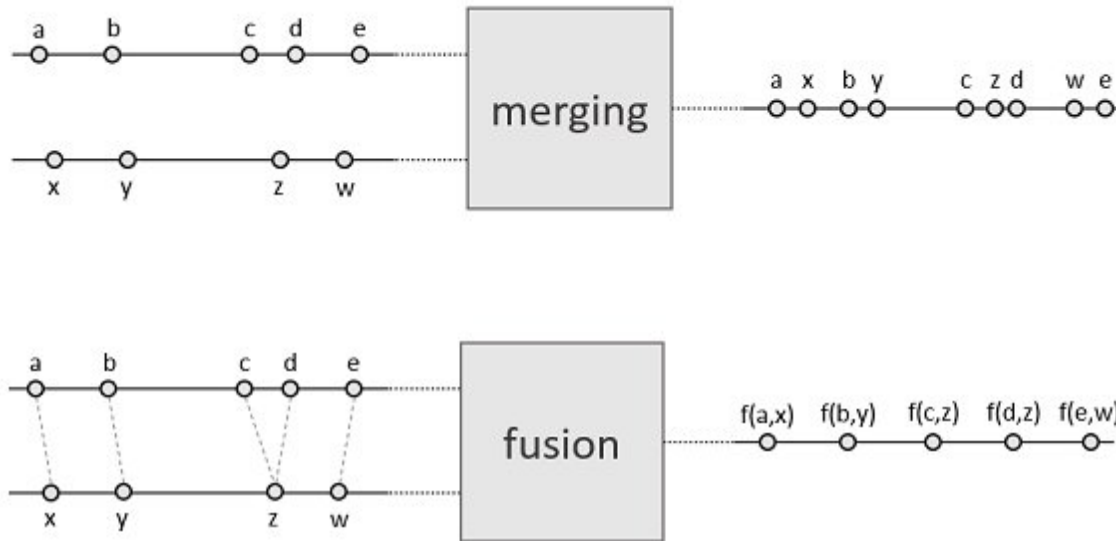
        // ToDo : Add more/modify stuff if needed

        // Constructor
        public Component(Pipeline pipeline, string name = "Component", DeliveryPolicy policy = null)
            : base(parent, name, policy)
        {
            // Connections creations
            InConnector = pipeline.CreateConnector<TYPE_IN>(nameof(InConnector));
            Out = pipeline.CreateEmitter<TYPE_OUT>(this, nameof(Out));

            // Registering method that will be called for every input
            InConnector.Out.Do(Process);
        }

        private void Process(TYPE_IN data, Envelope envelope)
        {
            //TO DO : process
            lock (this)
            {
                //If out:
                Out.Post(TYPE_OUT, DateTime.Now);
            }
        }
    }
}
```

► Deux classes de réceptions *complexes*



<https://github.com/microsoft/psi/wiki/Stream-Fusion-and-Merging>

- ▶ Deux opérateurs de réceptions *Merging*:
  - ▶ Zip: concaténation basé sur le temps dans les données.
  - ▶ Merge: concaténation rapide incluant les messages originaux (temps).
- ▶ Deux opérateurs de réceptions *Fusion*:
  - ▶ Fuse: opérateur générique
  - ▶ Join: opérateur spécialisé pour les interpolations 'reproductible'.
  - ▶ Interpolations:
    - Reproducible.Exact()
    - Reproducible.Nearest<T>(…)
    - Reproducible.First<T>(…)
    - Reproducible.Last<T>(…)
    - Reproducible.Linear(…)
    - Available.Exact()
    - Available.Nearest<T>(…)
    - Available.First<T>(…)
    - Available.Last<T>(…)
    - AdjacentValuesInterpolator(…)

- ▶ Il y a d'autres opérateurs :
  - ▶ Ceux de Linq (Select, First, Where, Flip, Count...)
  - ▶ Ceux concernant les temps (Delay, Parallel, Latency..)
  - ▶ Ceux concernant les streams (BridgeTo, PipeTo...)
  - ▶ Pair: concaténation des entrées basée sur le temps des messages du premier stream.
  - ▶ ...

## ► Liaison entre les composants :

```
var mic = new AudioCapture(p);
var resampler = new AudioResampler(p);
// link : micro -> sampler
mic.PipeTo(resampler);

var voiceActivityDetector = new SystemVoiceActivityDetector(p);
// link : sampler -> voiceDetector
resampler.PipeTo(voiceActivityDetector);

var azureSpeechReco = new AzureSpeechRecognizer(p);
// link : sampler + voiceDetector -> SpeechToText Azure
resampler.Join(voiceActivityDetector).PipeTo(azureSpeechReco);

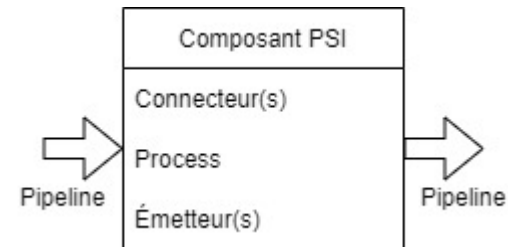
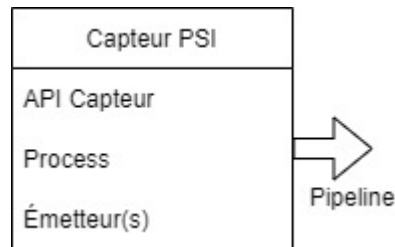
// Process : complete results only
var finalResults = azureSpeechReco.Out.Where(result => result.IsFinal);
NLPBot nlp = new NLPBot(p, botDescription);
// link : SpeechToText Azure -> NLP
finalResults.PipeTo(nlp);

var audioplayer = new AudioPlayer(p);
// link : NLP -> TextToSpeech
nlp.Voice.PipeTo(audioplayer);
```

- ▶ Exécution du pipeline:
  - ▶ Run(): exécute de façon bloquante le pipeline.
  - ▶ RunAsync(): non bloquant mais nécessite l'appel à *Dispose()*.

```
// RunAsync the pipeline in non-blocking mode.  
p.RunAsync(ReplayDescriptor.ReplayAllRealTime);  
// Waiting for an out key  
Console.WriteLine("Press any key to stop the application.");  
Console.ReadLine();  
// Stop correctly the pipeline.  
p.Dispose();
```

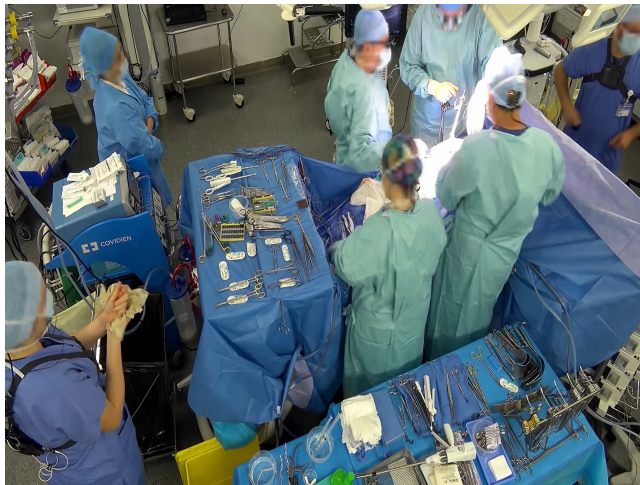
- ▶ Comme tout framework de haut niveau, il est facile de faire n'importe quoi.
- ▶ La courbe d'apprentissage se situe surtout sur l'optimisation des applications
- ▶ Deux 'types':
  - ▶ Les capteurs: émission
  - ▶ Les composants: réception & émission



- ▶ Toujours de l'émission -> Store !
- ▶ De préférence des tâches unitaires.

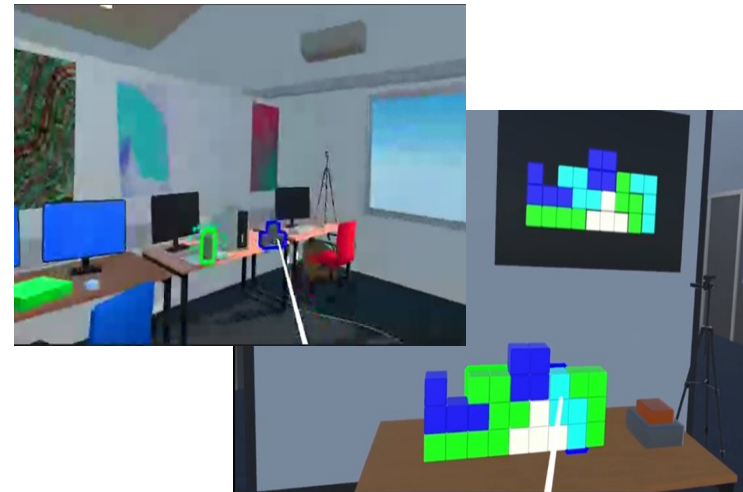


## Surgery



OR n°3 Centre Hépatobiliaire Paul Brousse

## VR task



A. Allemang--Trivalle  
ICMI 2023



A. Léchappé  
ICMI 2023



- **Critical, complex and stressful situation**
- **Team effort:** Surgeons, anesthesiologists, and nurses collaborate actively



© [ihorvsn] /Adobe Stock

Processes to integrate information about the situation of each team member

Use multimodal sensors to capture these processes (the situation) and reduce errors

- Collect precise behavioral data
- Enables fine control over the collaborative situation
- Deploy variety of environment augmentation in response to users' actions and collaborative states

## Machine-readable

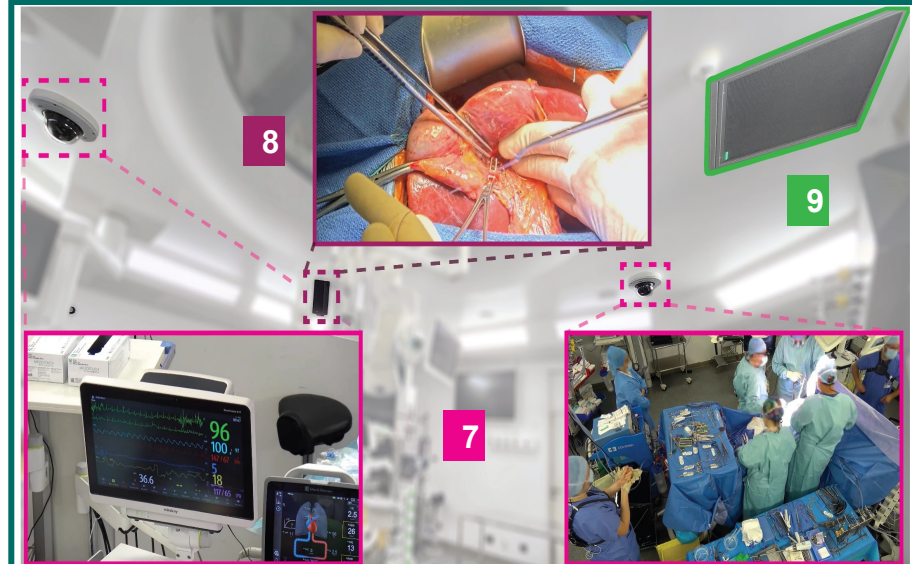
## Human-readable

### Wearable sensors



- 1** Lavalier microphone
- 2** Physiological garment
- 3** Tracking tag
- 4** Audio transmission system
- 5** Wi-Fi router
- 6** Computer (Microsoft PSI)

### Operating room sensors



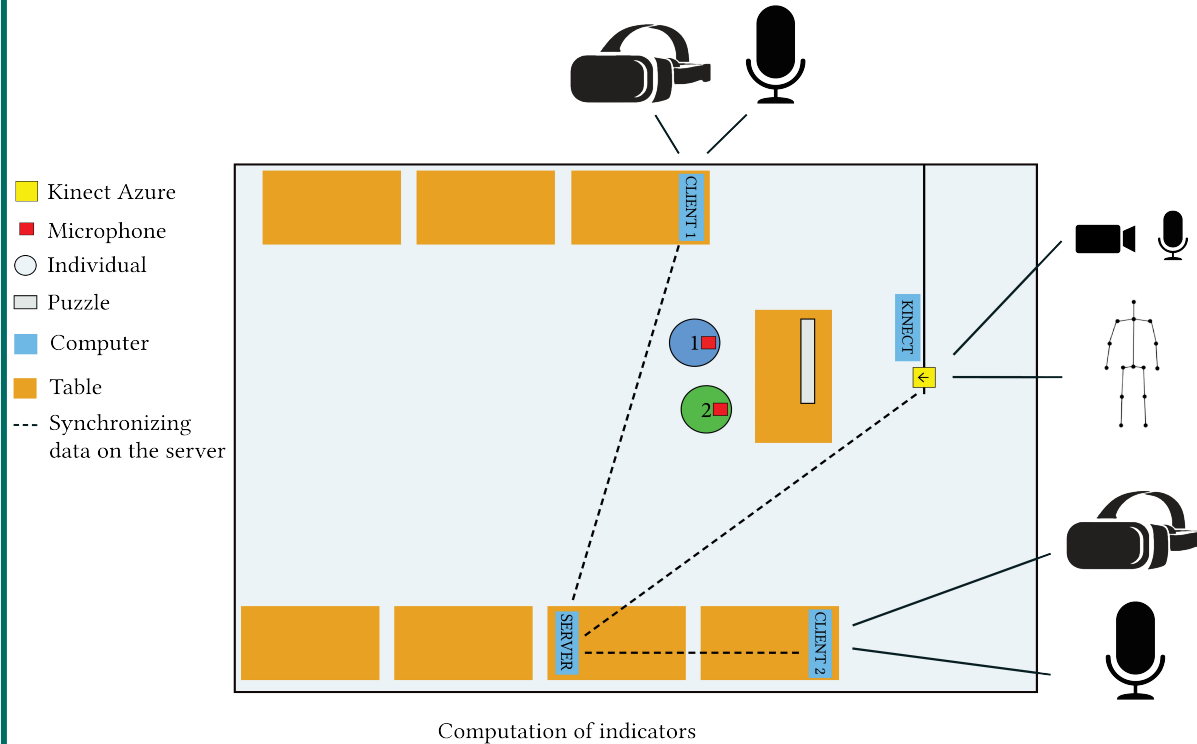
OR n°3 Centre Hépato-Biliaire Paul Brousse

- 7** PTZ cameras
- 8** Surgical camera
- 9** Ceiling microphone

## Collaborative activity

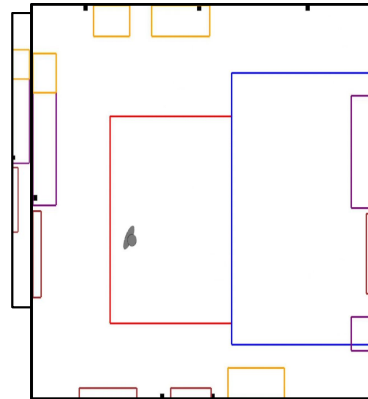
- Dyads
- Puzzle-solving task
  - Exploration/search task
  - Solving task
- Individuals must coordinate to perform tasks

## Sensors



## Metrics

- 1 Speech
- 2 Physiological data
- 3 Localization
- 4 Postures
- 5 Gaze



### User

- Speech flow
- Stress and mental workload
- Proximity to areas of interest
- Gaze on peers

### Team

- Overlapping Speech
- Turn-taking
- Dominance
- Proximity to team members
- Physiological synchrony
- Joint visual attention
- Mutual gaze

## Questionnaires

Get to know the people:

- **Background, Conditions, Relationships**

## Context

Annotations made by expert:

- **Task steps and phases**
- **Errors**
- **Adverse events** (surgery)
- **Collaboration quality**

# PSI - Exemple VR task

The screenshot displays a comprehensive VR task analysis interface. At the top, a timeline shows the session duration from 00:00:00 to 00:08:00, with a cursor at 00:05:00. Below the timeline are several views: an audio waveform, a video view of two participants in a VR lab, a 'Participants' skeletons view' showing two wireframe figures, and a 'Map view' of the virtual environment. The bottom section contains data tables for events, collaborative events, and interaction counts for both participants.

Event	Collaborative events	Example of data from Unity
User 2_Type : event_GAZED : SBlocExpe_Place	L'utilisateur 2 effectue GAZE sur ZBlocExpe_Place avec l'utilisateur 1	
Events on object	AvatarGaze	
User 2_GAZED : SBlocExpe_Place		
<b>Nombre d'interaction de P1 sur le puzzle</b>	<b>Nombre d'interaction de P2 sur le puzzle</b>	<b>Nombre de pièces collectées par P1</b>
Interactions : 6		
		<b>Nombre de pièces collectées par P2</b>
		Example of processed data

# UN PEU DE PRATIQUE...

```
// Create a pipeline
var p = Pipeline.Create();

// Create a timer component that produces a message every second
var timer = Timers.Timer(p, TimeSpan.FromSeconds(1));

// For each message created by the timer, print "Hello world!"
// along with the message's originating time.
timer.Do((t, e) =>
{
    Console.WriteLine($"{e.OriginatingTime}: Hello world!");
});

// Run the pipeline, but don't block here
p.RunAsync();

// Wait for the user to hit a key before closing the pipeline
Console.ReadKey();

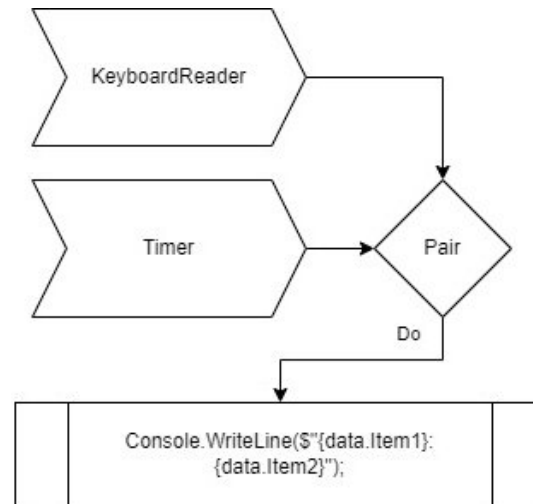
// Close the pipeline
p.Dispose();
```



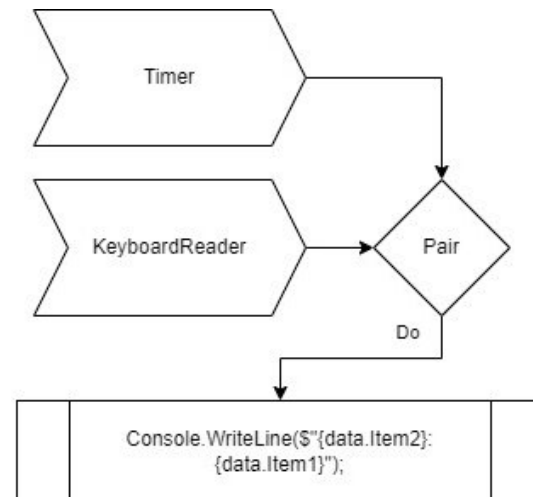


- ▶ Cloner <https://github.com/microsoft/psi-samples> et exécuter l'exemple [HelloWorld](#).
- ▶ Intégrer le composant *KeyboardReader* fournit ([lien](#)) au projet.
- ▶ Coder les fonctionnalités suivantes:
  - ▶ Afficher la sortie du *KeyboardReader* avec timestamp du message (on n'utilise pas le timer du projet pour le moment) :

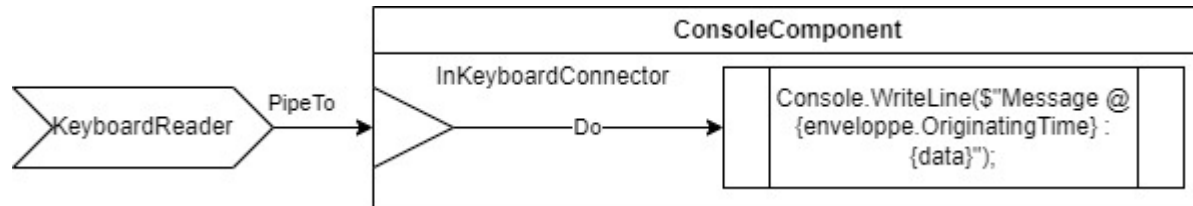
```
Console.WriteLine($"Message @ {envelope.OriginatingTime}: {data}")
```
  - ▶ Coder des fonctions qui prennent en entrée la sortie du *KeyboardReader* et la sortie du timer de l'exemple ([documentation](#)).
    - ▶ Une fonction doit afficher & répéter le dernier message du *KeyboardReader* à chaque *tick* du timer.



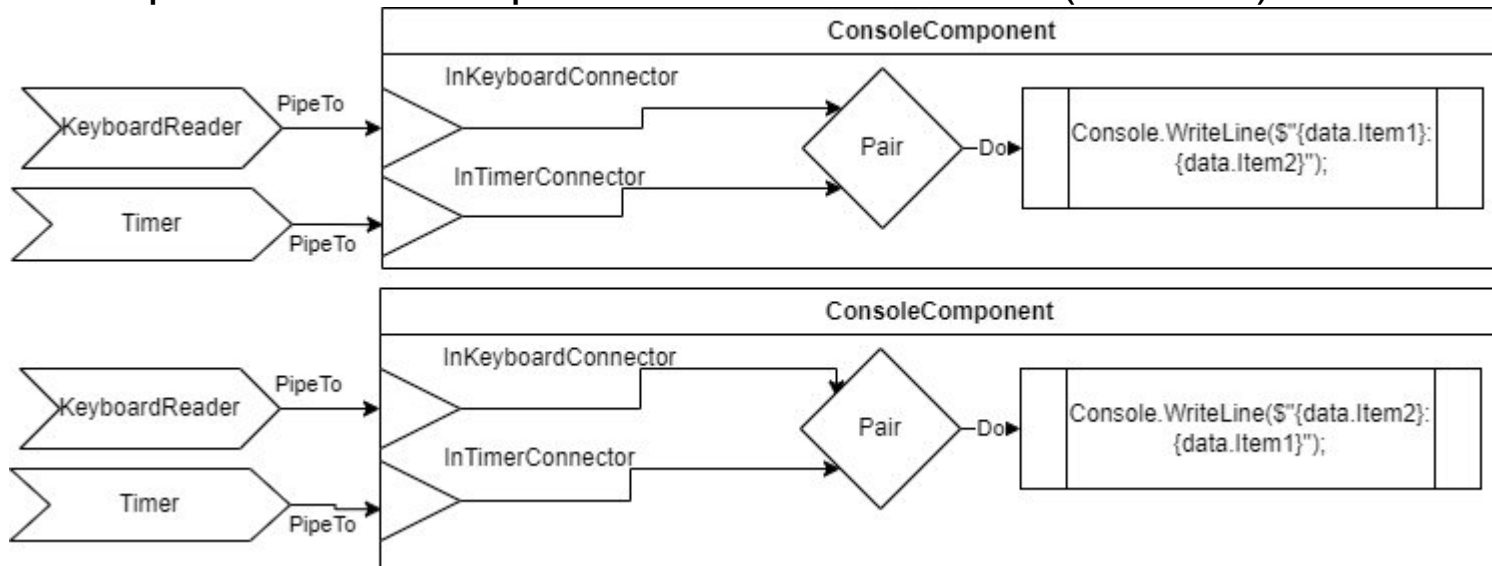
- L'autre doit afficher le dernier message du *KeyboardReader* avec le dernier *tick* du timer.



- ▶ Reprendre les fonctions précédemment réalisées sous formes de composants. *Un composant type est disponible dans les aides, [documentation](#).*

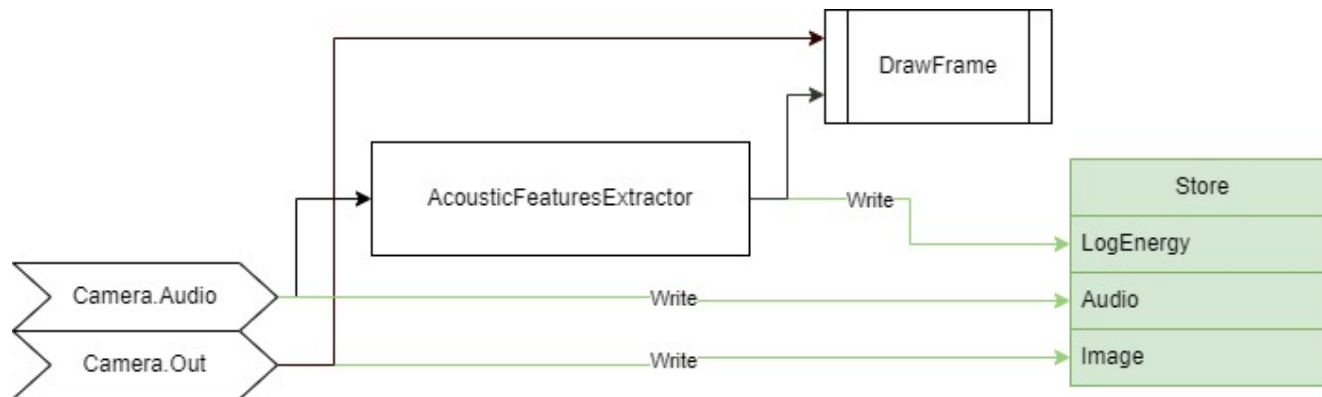


- ▶ Reprendre un exemple de fusion de données (au choix):



\psi

- ▶ Définir comme projet principal [WebcamWithAudioSample](#)
  - ▶ Exécuter et vérifier le bon fonctionnement.
  - ▶ Enregistrer les flux dans un(des) store(s) ([documentation](#)).



- ▶ Optimiser l'enregistrement avec la fonction [EncodeJpeg\(\)](#) :
  - ▶ Vérifiez que vous avez le package nuget dans le projet *Microsoft.Psi.Imaging.Windows*
  - ▶ Utiliser Psi Studio pour afficher en temps réel les informations stockées ([documentation](#)).

- ▶ Définir comme projet principal [AzureKinectSample](#)
- ▶ Intégrer l'enregistrement via un store de l'image couleur, de la profondeur et des squelettes.
- ▶ Afin d'afficher les squelettes, il est nécessaire d'ajouter le composant de visualisation de squelette (AzureKinect.Visualization.Windows.x64):

<https://github.com/microsoft/psi/wiki/3rd-Party-Visualizers>

- ▶ Celui-ci est disponible dans le dossier de PsiStudio fournit.
- ▶ Ouvrir dans Psi Studio et tester le mode *live*.

- ▶ De notre expérience, il est plus efficace d'effectuer un enregistrement par store, surtout quand il y a beaucoup de données (en fréquence et en taille).
- ▶ Effectuer un enregistrement où :
  - ▶ Deux personnes agissent en collaboration, on inclura des moments où les mains sont proches.
  - ▶ Un groupe de personne se rassemble autour d'un objet.
  - ▶ À votre imagination !

- ▶ Télécharger le projet : <https://github.com/AuMilliat/AzureKinectPostures>
- ▶ Développer une application :
  - ▶ Détection de groupes.
  - ▶ Détecter des positions d'intérêt : pointage, bras croisés, accroupie...
  - ▶ Détecter des actions de collaboration, comme le passage d'objet.
  - ▶ Lier ces détections à un enregistrement des images (couleur et profondeur).

- ▶ Les packages Unity sont disponible (Unity et WebRTC) :
  - ▶ <https://github.com/SaacPSI/saac>
- ▶ Les projets \psi fonctionnant avec Unity sont disponibles (PSIUnity & WebRTC) :
  - ▶ <https://github.com/AuMilliat/PsiTPCorrection>

