**CISCO SYSTEMS**

# Media Resource Control Protocol v2
# A Tutorial

**Sarvi Shanmugham,**

**Editor: MRCP v1/v2**

**Technical Leader, Cisco Systems**

1

# Roadmap

- **Overview of the IETF Speechsc WG Effort**

- **MRCP – Short Summary**

- **MRCP –Architecture Diagram**

- **MRCP - Usage**

- **MRCP v1 & v2 – Current Status**

# Overview of the IETF Speechsc WG Effort

- **IETF Working group - formed in 2002**

- **Aimed to develop a protocol that allows distributed speech processing(speech recognition, speaker recognition, verification and text-to-speech)**

- **Work with VoiceXML and SALT**

- **Leverage existing protocols as much as possible**

- **Leverage existing  W3C standards for markup**

# MRCP – Short Summary (contd.)

- **Basic Speech Services defined**

    **Speech Recognition**

    **Text-to-Speech**

    **Speaker Identification**

    **Speaker Verification**

    **Recording**

# MRCP – The Framework

- **The MRCP Framework leverages a suite of protocols and XML markup to achieve its purposes and only fills in where the needs have not already been addressed.**

    **SIP – This is used for discovering MRCP resources in the network and to rendezvous with the server and establish the necessary control and media pipes to the resources.**

    **SDP – SDP is used in conjunction with SIP for both resource discovery and the setup of control and media pipes for the session.**

    **RTP/RTCP – This is used for media transmission to/from the media processing resources.**

    **MRCP – This controls the operation of individual media processing resources, like ASR, TTS, SI, SV and recorders.**

# MRCP – The Framework (contd.)

- **W3C markup specifications**

    **SRGS – Definition of Voice Grammars that are processed by Speech Recognition engines.**

    **N-Grams – Stochastic Grammars.**

    **Semantic Tags – The above grammars could contain semantic markup associated with the grammars that aids in semantic processing of the recognized texts.**

    **SSML – Definitions Speech markup to be processed by Text-To-Speech Engines.**

    **NLSML – Natural Language Semantic Markup Language**
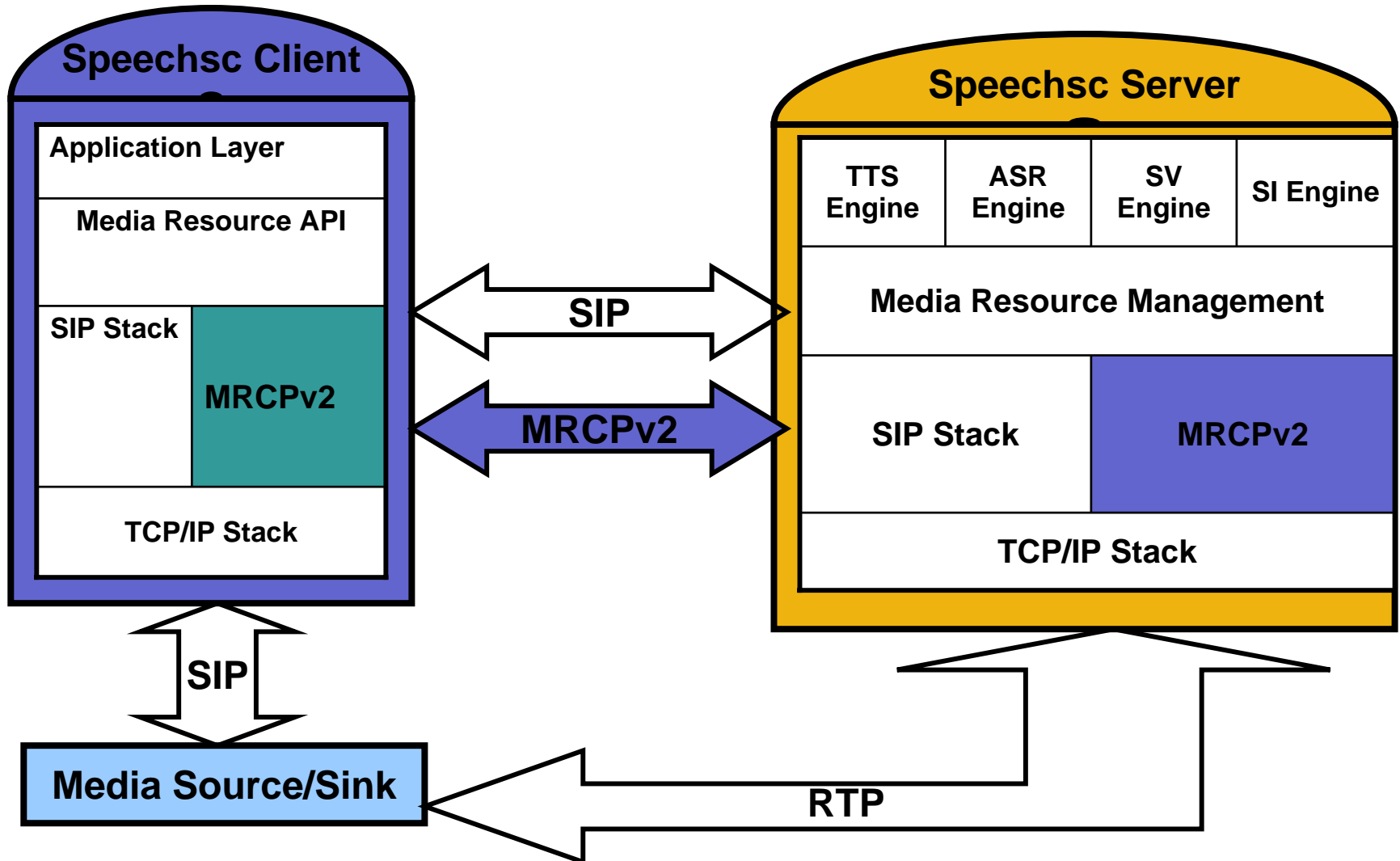
# MRCP – The Framework (contd.)

- **MRCP enhancements**

  **Recognition Results – The recognition resource returns results as a markup that is primarily based on NLSML. But there are a few minor additions to fill in gaps not addressed by NLML**

  **Grammar Enrollment Results – When enrolling new grammars, the results XML returned also contains extra information describing the enrollment status of the grammar enrollment.**

  **Speaker Identification/Verification Results – When doing Speaker Verification or Identification these XML extensions allow the resource to return the results of the verification or identification operation.**

# MRCP – Architecure Diagram

**Speechsc Client**

| Application Layer |
| Media Resource API |

| SIP Stack | MRCPv2 |

TCP/IP Stack

**Speechsc Server**

| TTS Engine | ASR Engine | SV Engine | SI Engine |

Media Resource Management

| SIP Stack | MRCPv2 |

TCP/IP Stack

SIP

MRCPv2

SIP

**Media Source/Sink**

RTP

# Server and Resource Addressing

- **Server**

  It's a regular SIP URI like the one below

  sip:mrcpv2@mediaserver.com

- **Resource Addressing**

  speechrecog - Speech Recognition

  dtmfrecog - DTMF Recognition

  speechsynth - Speech Synthesis

  basicsynth - Poorman's Speech Synthesizer

  speakverify - Speaker Verification

  recorder - Speech Recording

# MRCPv2 Protocol Basics

- **Connecting to the Server**

    **Uses a SIP INVITE and the SDP offer/answer model to connect to the media server and establish the session media and control pipes.**

    **Uses m= audio ….   For setting up media pipes to the server. This is the same as in any other SIP call setup.**

    **The m-line media stream established can shared by multiple mrcpv2 resource that may be part of the same SIP session.**

    **Uses m=control …. For setting up individual control pipes for each MRCPv2 resource that the client wants to control.**

    **There is one m=control .. line in the offer for every resource the client wants to allocate for the session.**

    **The m-lines specifies a transport type of TCP, SCTP or TLS and a fromat type of application/mrcpv2. The port number of this line MUST contain 9(discard port) in the offer and a valid server port in the answer. The client may then initiate an appropriate transport connection that port.**

# MRCPv2 Protocol Basics

- **Connecting to the Server**

    The offer m-line from the client also contains an "resource" specifying what type of resource it wants to allocate for the session. The corresponding answer m-line must contain a "channel" attribute that contains a channel identifier that will be used in all MRCP messages between the client and that specific resource.

    The transport connection(TCP, SCTP or TLS) could be shared across multiple MRCP sessions between a client and server.

- **Channel-Idenitifier**

    A channel identifier allocated for each resource is of the form

    > 32AECB234338@speechsynth

- **De-Allocating a Resource**

    To de-allocate a resource the client issues a SIP re-INVITE to the server where the appropriate m=control …. lines port is 0.

# MRCPv2 Protocol Basics

```
INVITE sip:mresources@mediaserver.com SIP/2.0
Via: SIP/2.0/TCP client.atlanta.example.com:5060;
     branch=z9hG4bK74bf9
Max-Forwards: 6
To: MediaServer <sip:mresources@mediaserver.com>
From: sarvi <sip:sarvi@cisco.com>;tag=1928301774
Call-ID: a84b4c76e66710
CSeq: 314161 INVITE
Contact: <sip:sarvi@cisco.com>
Content-Type: application/sdp
Content-Length: ...

v=0
o=sarvi 2890844526 2890842808 IN IP4 126.16.64.4
s=-
c=IN IP4 224.2.17.12
m=control 9 TCP application/mrcpv2
a=resource:speechsynth
a=cmid:1
m=audio 49170 RTP/AVP 0 96
a=rtpmap:0 pcmu/8000
a=recvonly
a=mid:1
```

# MRCPv2 Protocol Basics

```
SIP/2.0 200 OK
Via: SIP/2.0/TCP client.atlanta.example.com:5060;
        branch=z9hG4bK74bf9
To: MediaServer <sip:mresources@mediaserver.com>
From: sarvi <sip:sarvi@cisco.com>;tag=1928301774
Call-ID: a84b4c76e66710
CSeq: 314161 INVITE
Contact: <sip:sarvi@cisco.com>
Content-Type: application/sdp
Content-Length: ...

v=0
o=sarvi 2890844526 2890842808 IN IP4 126.16.64.4
s=-
c=IN IP4 224.2.17.12
m=control 32416 TCP application/mrcpv2
a=channel:32AECB234338@speechsynth
a=cmid:1
m=audio 48260 RTP/AVP 00 96
a=rtpmap:0 pcmu/8000
a=sendonly
a=mid:1
```

# MRCPv2 Protocol Basics

ACK sip:mresources@mediaserver.com SIP/2.0

Via: SIP/2.0/TCP client.atlanta.example.com:5060;

   branch=z9hG4bK74bf9 Max-Forwards: 6

To: MediaServer <sip:mresources@mediaserver.com>;tag=a6c85cf

From: Sarvi <sip:sarvi@cisco.com>;tag=1928301774

Call-ID: a84b4c76e66710

CSeq: 314162 ACK

Content-Length: 0

# Types of MRCP Messages

- **Request**

  **MRCP/2.0 434 SPEAK 543260**

  **Channel-Identifier: 32AECB23433802@speechsynth**

  **Voice-gender: neutral**

  **………**

- **Response**

  **MRCP/2.0 48 543260 200 IN-PROGRESS**

  **Channel-Identifier: 32AECB23433802@speechsynth**

  **………**

- **Event**

  **MRCP/2.0 73 SPEAK-COMPLETE 543260 COMPLETE**

  **Channel-Identifier: 32AECB23433802@speechsynth**

  **………**

# Generic Messages

- **Request**

    **SET-PARAMS**

    **GET-PARAMS**

- **Headers**

    **Channel-Identifier**

    **Active-Request-Id-List**

    **Proxy-Sync-Id**

    **Content-Id**

    **Content-Type**

    **Content-Length**

    **Content-Base**

    **Content-Location**

    **Content-Encoding**

    **Cache-Control**

    **Logging-Tag**

    **Set-Cookie**

    **Set-Cookie2**

    **Vendor-Specific**

# Text-To-Speech Resource

- **Request**

  **SPEAK**

  **STOP**

  **PAUSE**

  **RESUME**

  **BARGE-IN-OCCURRED**

  **CONTROL**

  **LOAD-LEXICON**

- **Event**

  **SPEECH-MARKER**

  **SPEAK-COMPLETE**

STOP

LOAD-LEXICON

**Idle**

SPEAK

STOP

SPEAK-COMPLETE

BARGE-IN-OCCURED

STOP

**Speaking**

CONTROL

RESUME

MARKER

PAUSE

**Paused**

CONTROL

PAUSE

# Text-To-Speech Resource

- **Headers**

| | |
|---|---|
| **Jump-Target** | **Fetch-hint** |
| **Kill-On-Barge-In** | **Audio-Fetch-Hint** |
| **Speaker-Profile** | **Fetch-Timeout** |
| **Completion-Cause** | **Failed-Uri** |
| **Completion-Reason** | **Failed-uri-cause** |
| **Voice-Parameter** | **Speak-Restart** |
| **Prosody-Parameter** | **Speak-Length** |
| **Speech-Marker** | **Load-Lexicon** |
| **Speech-Language** | **Lexicon-Search-Order** |

# Text-To-Speech Resource

## Speech Markup

```
<?xml version="1.0"?>

<speak>

<paragraph>

    <sentence> You have 4 new messages. </sentence>

    <sentence>The first is from <say-as type="name"> Stephanie
            Williams </say-as> and arrived at <break/>
            <say-as type="time"> 3:45pm </say-as>.

    </sentence>

    <sentence>The subject is <prosody rate="-20%"> ski
            trip </prosody>

    </sentence>

</paragraph>

</speak>
```
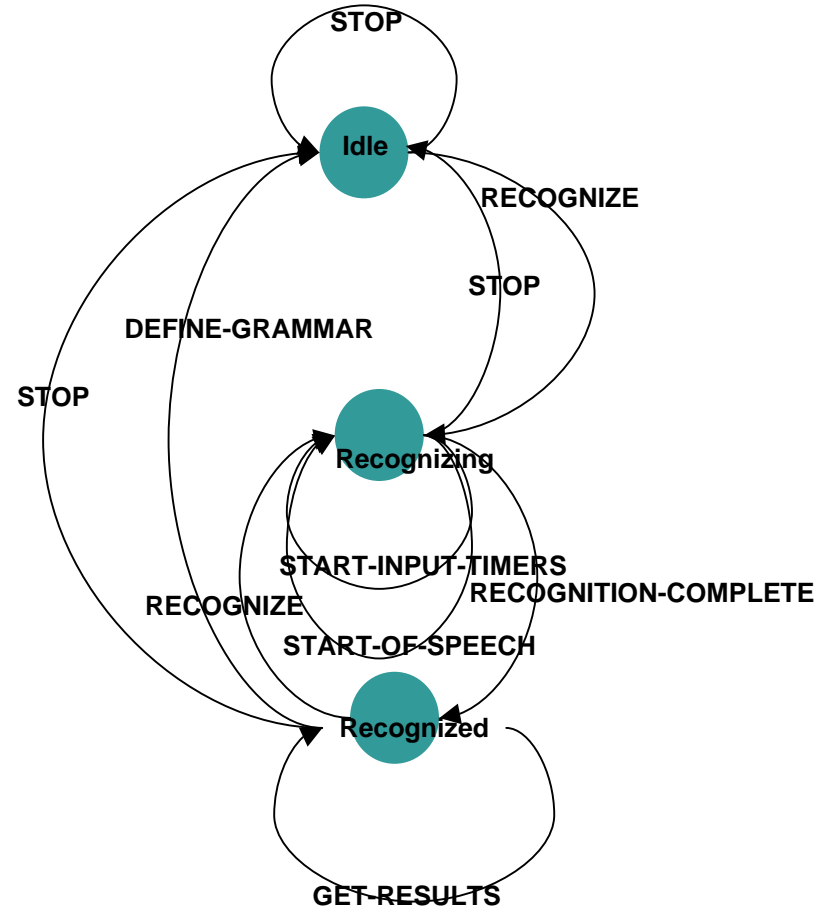
# Recognition Resource

- **Request**

    **DEFINE-GRAMMAR**

    **RECOGNIZE**

    **INTERPRET**

    **GET-RESULT**

    **START-INPUT-TIMERS**

    **STOP**

    **START-PHRASE-ENROLLMENT**

    **ENROLLMENT-ROLLBACK**

    **END-PHRASE-ENROLLMENT**

    **MODIFY-PHRASE**

    **DELETE-PHRASE**

- **Event**

    **START-OF-SPEECH**

    **RECOGNITION-COMPLETE**

    **INTERPRETATION-COMPLETE**

STOP

Idle

RECOGNIZE

STOP

DEFINE-GRAMMAR

STOP

Recognizing

START-INPUT-TIMERS

RECOGNITION-COMPLETE

RECOGNIZE

START-OF-SPEECH

Recognized

GET-RESULTS

# Recognition Resource

- **Recognition Headers**

  Confidence-Threshold

  Sensitivity-Level

  Speed-Vs-Accuracy

  N-Best-List-Length

  No-Input-Timeout

  Recognition-Timeout

  Waveform-Url

  Completion-Cause

  Completion-Reason

  Recognizer-Context-Block

  Start-Input-Timers

  Speech-Complete-Timeout

  Speech-Incomplete-Timeout

  Dtmf-Interdigit-Timeout

  Dtmf-Term-Timeout

  Dtmf-Term-Char

  Fetch-Timeout

  Failed-Uri

  Failed-Uri-Cause

  Save-Waveform

  New-Audio-Channel

  Speech-Language

  Ver-Buffer-Utterance

  Recognition-Mode

  Cancel-If-Queue

  Hotword-Max-Duration

  Hotword-Min-Duration

  Interpret-text

# Recognition Resource

- **Enrollment Headers**

  **Num-Min-Consistent-Pronunciations**

  **Consistency-Threshold**

  **Clash-threshold**

  **Personal-Grammar-Uri**

  **Phrase-Id**

  **Phrase-NL**

  **Weight**

  **Save-Best-Waveform**

  **New-Phrase-Id**

  **Confusable-Phrases-Uri**

  **Abort-Phrase-Enrollment**

# Recognition Resource

## Grammar Markup

```
<?xml version="1.0"?>

<!-- the default grammar language is US
  English -->

<grammar xml:lang="en-US" version="1.0">

<!-- single language attachment to
  tokens -->

<rule id="yes">

           <one-of>

                <item xml:lang="fr-
  CA">oui</item>

                <item xml:lang="en-
US">yes</item>

           </one-of>

</rule>

<!-- single language attachment to a
  rule expansion -->

<rule id="request">

           may I speak to

           <one-of xml:lang="fr-CA">

                <item>Michel
  Tremblay</item>

                <item>Andre Roy</item>

           </one-of>

</rule>
```

```
<!-- multiple language attachment to a
  token -->

<rule id="people1">

           <token lexicon="en-US,fr-
CA"> Robert </token>

</rule>

<!-- the equivalent single-language
attachment

           expansion -->

<rule id="people2">

           <one-of>

                <item xml:lang="en-
US">Robert</item>

                <item xml:lang="fr-
CA">Robert</item>

           </one-of>

</rule>

</grammar>
```

# Recognition Resource

## Result Markup

```xml
<?xml version="1.0"?>
<result
 grammar="http://theYesNoGrammar">
    <interpretation>
        <instance>
            <myApp:yes_no>

  <response>yes</response>
            </myApp:yes_no>
        </instance>
        <input>ok</input>
    </interpretation>
</result>
```

# Recognition Resource

## Enrollment Result Markup

```
<?xml version= "1.0"?>

<result grammar="Personal-Grammar-URI"
  xmlns:mrcp=
        "http://www.ietf.org/mrcp2">

<mrcp:result-type type="ENROLLMENT"/>
 <mrcp:enrollment-result>
      <num-clashes> 2 </num-clashes>
      <num-good-repetitions> 1
          </num-good-repetitions>
      <num-repetitions-still-needed> 1
          </num-repetitions-still-needed>
      <consistency-status> consistent
          </consistency-status>
      <clash-phrase-ids>
    <item> Jeff </item>
    <item> Andre </item>
</clash-phrase-ids>
```

```
<transcriptions>
     <item> m ay b r ow k er </item>
     <item> m ax r aa k ah </item>
   </transcriptions>
   <confusable-phrases>
     <item>
        <phrase> call </phrase>
        <confusion-level> 10
            </confusion-level>
     </item>
   </confusable-phrases>
</mrcp:enrollment-result>
</result>
```

# Recording Resource

- **Request**
  - **RECORD**
  - **STOP**
  - **START-INPUT-TIMERS**

- **Event**
  - **START-OF-SPEECH**
  - **RECORD-COMPLETE**

- **Headers**

  ```
  Sensitivity-Level        Media-Type

  No-Input-Timeout         Max-Time

  Completion-Cause         Final-Silence

  Completion-Reason        Capture-On-Speech

  Failed-Uri               Ver-Buffer-Utterance

  Failed-Uri-Cause         Start-input-timers

  Record-Uri               New-audio-channel
  ```

STOP

Idle

STOP

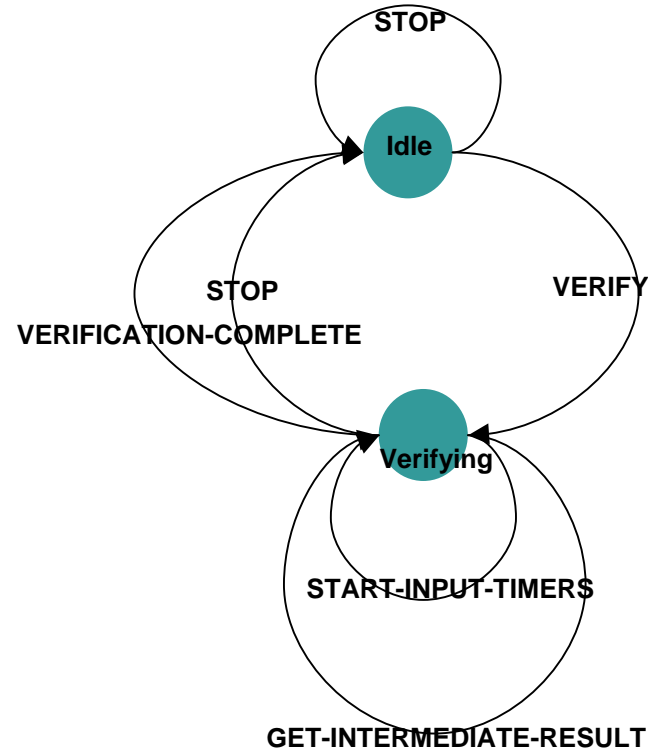RECORD-COMPLETE

RECORD

Recording

# Verification Resource

- **Request**

    START-SESSION

    END-SESSION

    QUERY-VOICEPRINT

    DELETE-VOICEPRINT

    VERIFY

    VERIFY-FROM-BUFFER

    VERIFY-ROLLBACK

    STOP

    CLEAR-BUFFER

    START-INPUT-TIMERS

    GET-INTERMEDIATE-RESULT

- **Event**

    VERIFICATION-COMPLETE

    START-OF-SPEECH

# Verification Resource

- **Verification Headers**

`Repository-Uri`

`Voiceprint-Identifier`

`Verification-Mode`

`Adapt-Model`

`Abort-Model`

`Security-Level`

`Num-Min-Verification-Phrases`

`Num-Max-Verification-Phrases`

`No-Input-Timeout`

`Save-Waveform`

`Waveform-Url`

`Voiceprint-Exists`

`Ver-Buffer-Utterance`

`Input-Waveform-Url`

`Verification-Type`

`Digit-Sequence`

`Completion-Cause`

`Completion-Reason`

`Speech-Complete-Timeout`

`New-Audio-Channel`

`Abort-Verification`

**Start-Input-Timers**

# Verification Resource

## Verification Result Markup

```
<?xml version="1.0"?>
<result grammar="What-Grammar-URI"
  xmlns:mrcp="http://www.ietf.org/mrcp2">
  <mrcp:result-type type="VERIFICATION" />
  <mrcp:verification-result>
    <voiceprint id="johnsmith">
      <adapted> true </adapted>
      <incremental>
        <num-frames> 50 </num-frames>
        <device> cellular-phone </device>
        <gender> female </gender>
        <decision> accepted </decision>
        <verification-score> 0.98514 </verification-score>
      </incremental>
      <cumulative>
        <num-frames> 1000 </num-frames>
        <device> cellular-phone </device>
```

# Verification Resource

## Verification Result Markup(contd.)

```
        <gender> female </gender>

        <decision> accepted </decision>

        <verification-score> 0.91725</verification-score>

      </cumulative>

    </voiceprint>

    <voiceprint id="marysmith">

      <cumulative>

        <verification-score> 0.93410 </verification-score>

      </cumulative>

    </voiceprint>

    <voiceprint uri="juniorsmith">

      <cumulative>

        <verification-score> 0.74209 </verification-score>

      </cumulative>

    </voiceprint>

  </mrcp:verification-result>

</result>
```

# Call Flow Example

```
C->S:

INVITE sip:mresources@mediaserver.com SIP/2.0

Max-Forwards: 6

To: MediaServer <sip:mresources@mediaserver.com>

From: sarvi <sip:sarvi@cisco.com>;tag=1928301774

Call-ID: a84b4c76e66710

CSeq: 314163 INVITE

Contact: <sip: sarvi@cisco.com>

Content-Type: application/sdp

Content-Length: 142


v=0

o=sarvi 2890844526 2890842809 IN IP4 126.16.64.4

s=SDP Seminar

i=A session for processing media

c=IN IP4 224.2.17.12/127

m=control 9 SCTP application/mrcpv2

a=resource:speechsynth

a=cmid:1

m=audio 49170 RTP/AVP 0 96

a=rtpmap:0 pcmu/8000

a=recvonly

a=mid:1
```

```
m=control 9 SCTP application/mrcpv2

a=resource:speechrecog

a=cmid:2

m=audio 49180 RTP/AVP 0 96

a=rtpmap:0 pcmu/8000

a=rtpmap:96 telephone-event/8000

a=fmtp:96 0-15

a=sendonly

a=mid:2
```

# Call Flow Example

```
S->C:

SIP/2.0 200 OK

To: MediaServer <sip:mresources@mediaserver.com>

From: sarvi <sip:sarvi@cisco.com>;tag=1928301774

Call-ID: a84b4c76e66710

CSeq: 314163 INVITE

Contact: <sip: sarvi@cisco.com>

Content-Type: application/sdp

Content-Length: 131


v=0

o=sarvi 2890844526 2890842809 IN IP4 126.16.64.4

s=SDP Seminar

i=A session for processing media

c=IN IP4 224.2.17.12/127

m=control 32416 SCTP application/mrcpv2

a=channel:32AECB23433801@speechsynth

a=cmid:1

m=audio 48260 RTP/AVP 0

a=rtpmap:0 pcmu/8000

a=sendonly

a=mid:1
```

```
m=control 32416 SCTP application/mrcpv2

a=channel:32AECB23433802@speechrecog

a=cmid:2

m=audio 48260 RTP/AVP 0

a=rtpmap:0 pcmu/8000

a=rtpmap:96 telephone-event/8000

a=fmtp:96 0-15

a=recvonly

a=mid:2


C->S:

ACK sip:mrcp@mediaserver.com SIP/2.0

Max-Forwards: 6

To: MediaServer
      <sip:mrcp@mediaserver.com>;tag=a6c85cf

From: Sarvi <sip:sarvi@cisco.com>;tag=1928301774

Call-ID: a84b4c76e66710

CSeq: 314164 ACK

Content-Length: 0
```

# Call Flow Example

```
C->S: MRCP/2.0 386 SPEAK 543257

Channel-Identifier: 32AECB23433802@speechsynth

Kill-On-Barge-In: false

Voice-gender: neutral

Voice-category: teenager

        Prosody-volume: medium

Content-Type: application/synthesis+ssml

Content-Length: 104


<?xml version="1.0"?>

<speak>

<paragraph>

        <sentence>You have 4 new
    messages.</sentence>

        <sentence>The first is from <say-as

        type="name">Stephanie Williams</say-as>
    <mark name="Stephanie"/>

        and arrived at <break/>

        <say-as type="time">3:45pm</say-
    as>.</sentence>


        <sentence>The subject is <prosody

        rate="-20%">ski
    trip</prosody></sentence>

</paragraph>

</speak>
```

```
S->C: MRCP/2.0 49 543257 200 IN-PROGRESS

Channel-Identifier: 32AECB23433802@speechsynth


S->C: MRCP/2.0 46 SPEECH-MARKER 543257 IN-
    PROGRESS

Channel-Identifier: 32AECB23433802@speechsynth

        Speech-Marker: Stephanie


The synthesizer finishes with the SPEAK request.


S->C: MRCP/2.0 48 SPEAK-COMPLETE 543257 COMPLETE

Channel-Identifier: 32AECB23433802@speechsynth
```

# Call Flow Example

```
C->S:MRCP/2.0 343 RECOGNIZE 543258

Channel-Identifier: 32AECB23433801@speechrecog

Content-Type: application/grammar+xml

Content-Length: 104


<?xml version="1.0"?>


<!-- the default grammar language is US English
     -->
<grammar xml:lang="en-US" version="1.0">


<!-- single language attachment to a rule
     expansion -->
     <rule id="request">
          Can I speak to
          <one-of xml:lang="fr-CA">
                    <item>Michel Tremblay</item>

                    <item>Andre Roy</item>

          </one-of>
     </rule>


</grammar>


S->C: MRCP/2.0 49 543258 200 IN-PROGRESS
Channel-Identifier: 32AECB23433801@speechrecog
```

```
C->S: MRCP/2.0 289 SPEAK 543259

Channel-Identifier: 32AECB23433802@speechsynth

          Kill-On-Barge-In: true

Content-Type: application/sml

Content-Length: 104


<?xml version="1.0"?>
<speak>
<paragraph>
          <sentence>Welcome to ABC
     corporation.</sentence>
          <sentence>Who would you like Talk
     to.</sentence>
</paragraph>
</speak>


S->C: MRCP/2.0 52 543259 200 IN-PROGRESS
Channel-Identifier: 32AECB23433802@speechsynth
```

# Call Flow Example

```
S->C: MRCP/2.0 49 START-OF-SPEECH 543258 IN-PROGRESS

Channel-Identifier: 32AECB23433801@speechrecog

        Proxy-Sync-Id: 987654321


C->S: MRCP/2.0 69 BARGE-IN-OCCURRED 543259

Channel-Identifier: 32AECB23433802@speechsynth

        Proxy-Sync-Id: 987654321


S->C: MRCP/2.0 72 543259 200 COMPLETE

Channel-Identifier: 32AECB23433802@speechsynth

        Active-Request-Id-List: 543258


S->C: MRCP/2.0 73 SPEAK-COMPLETE 543259 COMPLETE

Channel-Identifier: 32AECB23433802@speechsynth

        Completion-Cause: 001 barge-in


S->C: MRCP/2.0 412 RECOGNITION-COMPLETE 543258 COMPLETE

Channel-Identifier: 32AECB23433801@speechrecog

        Completion-Cause: 000 success

Waveform-URL: http://web.media.com/session123/audio.wav

Content-Type: application/x-nlsml

Content-Length: 104
```
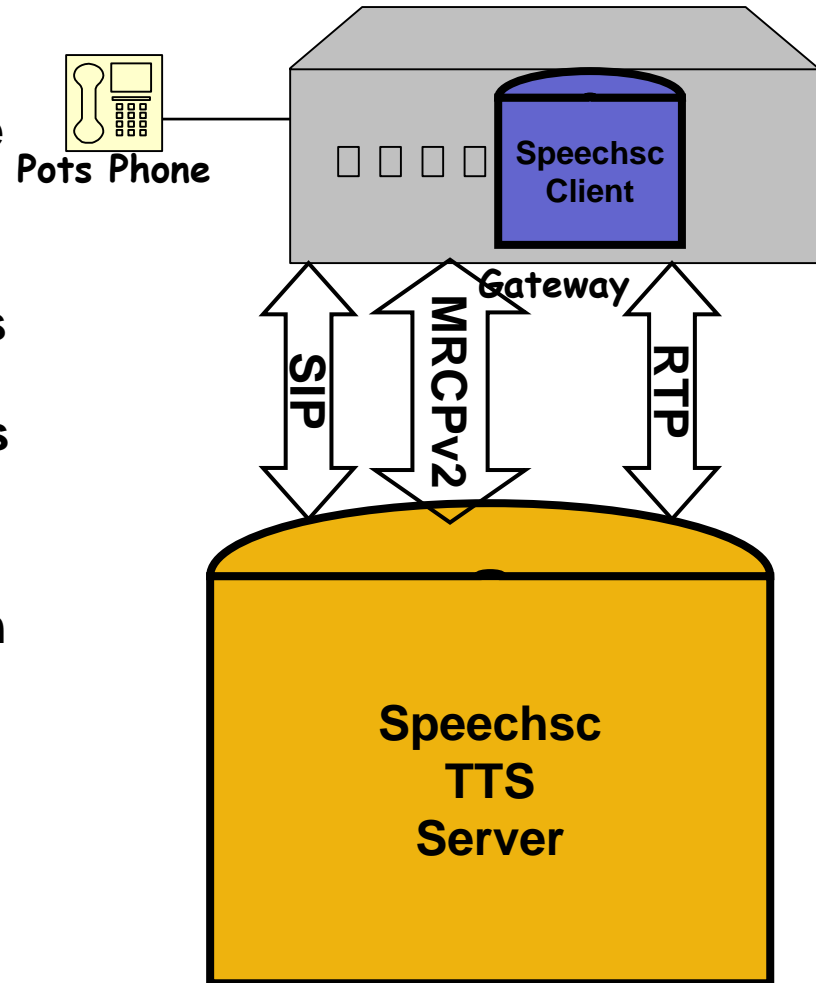
```
<?xml version="1.0"?>
<result x-model="http://IdentityModel"
  xmlns:xf="http://www.w3.org/2000/xforms"
  grammar="session:request1@form-level.store">
    <interpretation>
        <xf:instance name="Person">
            <Person>
                <Name> Andre Roy </Name>
            </Person>
        </xf:instance>
        <input>   may I speak to Andre Roy </input>
    </interpretation>
</result>


C->S:BYE sip:mrcp@mediaserver.com SIP/2.0

Max-Forwards: 6

From: Sarvi <sip:sarvi@cisco.com>;tag=a6c85cf

To: MediaServer
      <sip:mrcp@mediaserver.com>;tag=1928301774

Call-ID: a84b4c76e66710

CSeq: 231 BYE

Content-Length: 0
```
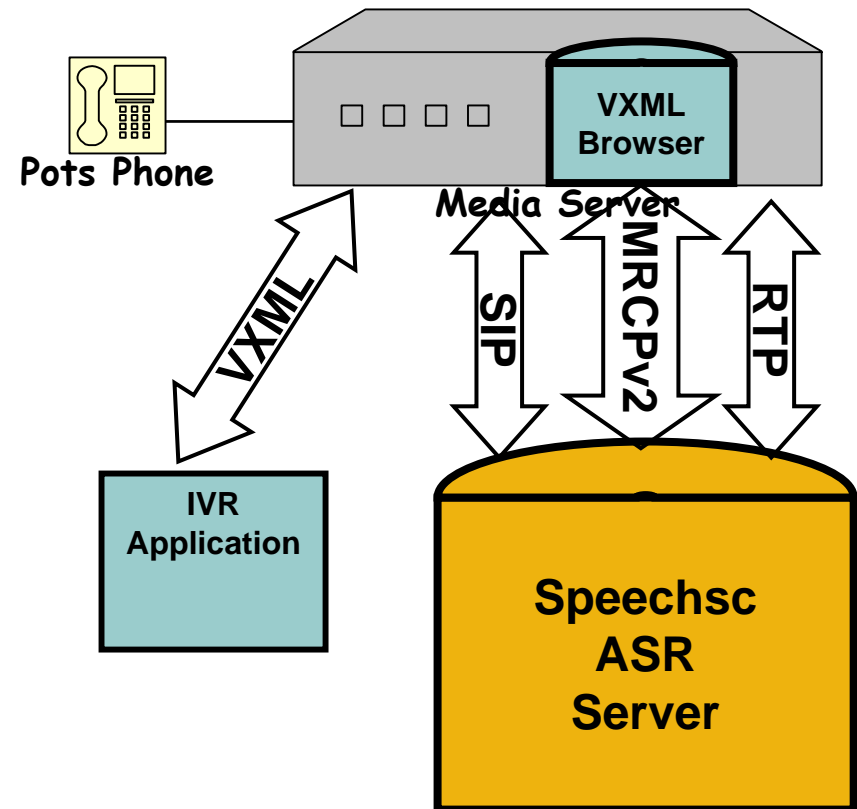
# Use Case: Text to Speech Announcements

- **POTS phone attempts call.**

- **VoIP gateway, acting as a SIP UA, attempts SIP session to complete the call; gets error, like "486 Busy Here".**

- **VoIP Gateway constructs a text error string from the SIP message, such as "Your call to 978-555-1212 did not go through because the called party was busy".**

- **Gateway INVITES SPEECHSC server to connect RTP stream and issues an MRCPv2 TTS request for the error message**

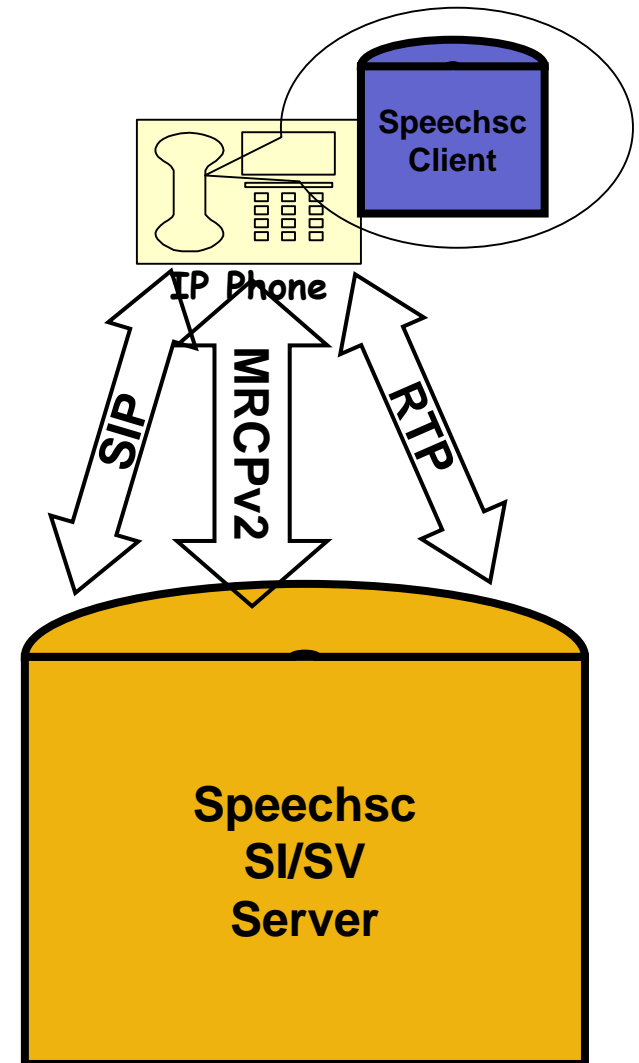- **Speechsc server plays message to the user on the POTS phone.**

**Pots Phone**

**Speechsc Client**

**Gateway**

**SIP**

**MRCPv2**

**RTP**

**Speechsc TTS Server**

# Use Case: VXML-based ASR

- **Users call into the service in order to obtain stock quotes.**

- **Media Server fetches VoiceXML to drive user interaction.**

- **Media Server INVITEs Speechsc server for ASR**

- **VoiceXML interpreter on the Media Server directs the user's media stream to the ASR server and uses MRCPv2 to control the ASR server.**

- **Results come back and the application proceeds.**

Pots Phone

VXML Browser

Media Server

VXML

SIP

MRCPv2

RTP

IVR Application

Speechsc ASR Server

# Use Case: Speaker Verification

- **A user speaks into a SIP phone to "log in" to that phone to make and receive phone calls using his identity and preferences**

- **IP phone uses SIP and MRCPv2 to set up an RTP stream between the phone and the SPEECHSC SI/SV server and request verification.**

- **SV server verifies the user's identity and returns the result via MRCPv2.**

- **The IP Phone may either use the identity directly to identify the user in outgoing calls, to fetch the user's preferences from a configuration server, request authorization from a AAA server, etc.**

**Speechsc Client**

**IP Phone**

**SIP**

**MRCPv2**

**RTP**

**Speechsc SI/SV Server**

# Current WG Status

- **Requirements Document passed IESG Review - soon to be published as an RFC**

  **draft-ietf-speechsc-reqts-05.txt**

- **MRCPv2 Protocol Document in second revision - expect last call in late fall**

  **draft-ietf-speechsc-mrcpv2-04.txt**

- **MRCPv1 Protocol Document is pending IESG review for publication as an Informational RFC.**

  **http://www.ietf.org/internet-drafts/draft-shanmugham-mrcp-05.txt**