PositionPaperfortheW3CMMIWorkshop
GrahamWilcockandKristiinaJokinen
UniversityofHelsinkiandUniversityofTampere

SCXML,MultimodalDialogueSystemsandMMIArchitecture

Weareinterestedintheworkshopbecausethetopicisimportantforour         research:we
wishtolearnmoreabouthowtheMMIarchitecturesupports
- a) fusionofmodalities
- b) incrementalpresentation
- c) designofcooperativeinteraction.

ThefirstauthorhasworkedonXML       -basedlanguageprocessing(naturallanguage
analysis,generation,  andannotation)andisinterestedinpracticalapplicationsof
emergingtechnology,whilethesecondauthorhasplayedaleadingroleinseveral
researchanddevelopmentprojectsonspokendialoguesystems,andisinterestedin
prinicpled-basedrepresentat ionandarchitecturesforimplementingfundamental
aspectsofhumancommunication.Ourcurrentpurposeistocombineourprevious
work.WeapproachMMIfromtwodifferentbutrelatedperspectives:SCXMLasa
basisforvoiceinterfacesandcooperativemulti       modalroutenavigation.

SCXMLasabasisforvoiceinterfacesisdescribedin[Wilcock2007].Thebasicideas
ofstatechartsareillustratedbymeansofasimple"stopwatch"demowithanopen
sourceJavaimplementationofSCXML.Thebasicversionofthe        democomesfrom
JakartaCommonsSCXML[ApacheSoftwareFoundation2006],anopensource
JavaimplementationofSCXML,andhasagraphicaluserinterface(GUI),which
displaysthetimeandenablesthestopwatchtobestarted,paused,stoppedandreset
bymo useclicks.Thevoiceuserinterface(VUI)isaddedtothedemousingthe
Sphinx-4opensourceJavaspeechrecognizer[CarnegieMellonUniversity2004]and
theFreeTTSopensourceJavaspeechsynthesizer[SunMicrosystems2005].Ina
"speakingstopwatch"v ersionofthedemo,whentheuserstopsorpausesthe
stopwatchbyamouseclickontheappropriatebutton,thetimeisreadoutaloudby
thespeechsynthesizer.Inaddition,abriefpromptisspokenwhentheuserstarts,
un-pausesorresetsthestopwatch.    Ina"listeningstopwatch"versionofthedemo,
theuserstarts,stops,pauses,un     -pauses,andresetsthestopwatcheitherbyvoice
commandsorbymouseclicks.ThespeechrecognizerusesasmallJSGFgrammar
forthevoicecommands.Allthreeversions(ba        sic,speaking,listening)followthe
samestatetransitionswhicharedefinedintheSCXMLstatechartfile.

CooperativemultimodalroutenavigationisthebasisoftheMUMSsystem,aPDA         -
basedroutenavigationsystemwhichallowstheusertoquerypublic         transportation
informationusingspokenlanguagecommandsandpen        -pointinggesturesonamap.
Italsoprovidesrouteinformationinspeechandgraphicaloutput.Thesystemis
describedinmoredetailin[Hurtig&Jokinen2005,2006;Jokinen2007],andits
evaluationisreportedin[Jokinen&Hurtig2006].However,althoughtheMUMS
systemisbasedontheSOA(ServiceOrientedArchitectureapproach,andXMLis
usedasageneralinterfacelangaugebetweenthedifferentlayersoftechnology
components,itis   stillratherapplicationspecificinitsrepresentationandprocessing
ofinformation.Wearethuslookingforamoresystematicapproachtodevelop

applications, so that the infrastructure would allow easy experimentation with different technology compone nts and their internal functioning.

Since SCXML supports a clean separation of data, logic, and user interface, based on the data-flow-presentation (DFP) architectural pattern, we believe it has benefits for our research. Our ultimate goal is to experime nt with different possibilities to develop natural user interfaces; natural in the sense of allowing the use of natural language, and also in the sense of providing intuitive functionality for the user. Also, we are interested in investigating how the desi gn of spoken dialogue systems and multimodal route navigation system with an emphasis on the human cooperation aspects will be enabled in SCXML and MMI. In dialogue research, SCXML has been proposed by [Kronlid and Lager 2007] as a basis for implementing t he information-state update (ISU) framework for dialogue management, which is a promising approach we would like to pursue further.

The route navigation application resembles the "DrivingDirections" user case of W3C Multimodal Interaction Use Cases sinc e in both cases the system gives instructions of how to go forward. In both cases, the user needs to understand the instructions given by the system, and the system should "listen" to the user and observe if the message has gone across.

We are especially keen on finding good solutions for the interaction problems that are generally considered "natural" (in the above two senses) but which are also generally considered difficult due to lack of descriptive research and available technology:
(1) incremental repr esentation of information and allowing the user to zoom in and out both verbally and on the map
(2) allowing users to give feedback concerning their understanding in different ways: providing an answer to an explicit question ("Did you say the Opera stop?"), continuing the interaction with an appropriate next step ("Give me the next piece of information"), and by subtle signalling in their speech (variation of pronunciation together with the length of the following pause can signal wish to continue rather than the end of one's turn).

Concerning the topics of the workshop listed in the CFP, we would in particular like to address the following questions:

- Requirements for extensions to the MMI Architecture to improve the support of speech, GUI and Ink interface so n portable handheld multimodal devices.
- How to dynamically select appropriate modalities.
- Use of scripts to enable the customization of the user interface based upon previous user input.
- Support for effective user interfaces for various modes of interactio n, in terms of contextual prompts, constrained text input, and declarative event handlers, taking account of uncertainties in user input.

We are also interested in the solution to the following questions:

- How to process early and late information fusion.
- Plans to support multimodal applications and what standards are needed.
- Re-use of existing markup languages for prompts and constraints on user input.

References

ApacheSoftwareFoundation:TheApacheJakartaProject,CommonsSCXML. http://jakarta.apache.org/commons/scxml/(2006)

CarnegieMellonUniversity:Sphinx  -4:Aspeechrecognizerwrittenentirelyin theJavaprogramminglanguage.http://cmusphinx.sourceforge.net/sphinx4/(2004)

Harel,D.:Statecharts:AVisualFormalismforComplexSystems.Sc      ienceof ComputerProgramming8,NorthHolland(1987).

Hurtig,T.andJokinen,K.:OnMultimodalRouteNavig       ationinPDAs, $2^{nd}$Baltic ConferenceonH  umanLanguageTechnologies,pp.261    –266, Tallinn,Estonia,2005.

Hurtig,T.andJokinen,K.:    Modalityfus ioninaroutenavigationsystem.Proceedings oftheIUI2006WorkshoponEffectiveMultimodaDialogueIInterfaces,pp.19         -24, 2006.

Jokinen,K.:InteractionandMobileRouteNavigationApplication.      InMeng,L.,A.Zipf, andS.Winter(eds.)   Map-basedmob ileservices  -usagecontext,interactionand application,SpringerseriesonGeoinformatics,2007.

Jokinen,K.andHurtig,T   .:UserExpectationsandRealExperienceonaMultimodal InteractiveSystem.ProceedingsofInterspeech,Pittsburgh,US,2006.

Kronlid,F.andLager,T.:    ImplementingtheInformation -StateUpdateApproachto DialogueManagementinaSlightlyExtendedSCXML.    11thInternationalWorkshop ontheSemanticsandPragmaticsofDialogue     ,Trento,Italy.pp.99  -106,2007.

SunMicrosystems:Fr  eeTTS1.2:Aspeechsynthesizerwrittenentirelyinthe Javaprogramminglanguage.http://freetts.sourceforge.net/(2005)

Wilcock,G.:SCXMLandVoiceInterfaces.3    [rd]BalticConferenceonHuman LanguageTechnologies,Kaunas,Lithuania,2007. http://conference.vdu.lt/viewabstract.php?id=158&cf=7

W3C:StateChartXML(SCXML):StateMachineNotationforControlAbstraction. http://www.w3.org/TR/2007/WD-scxml-20070221/(2007)