

Kommunikative Rhythmen in Gestik und Sprache

Ipke Wachsmuth

Technische Fakultät, Universität Bielefeld
D-33594 Bielefeld
ipke@techfak.uni-bielefeld.de

Abstract. Gestik und Sprache sind die Eckpfeiler in der natürlichen Verständigung zwischen Menschen. Nicht von ungefähr wird ihnen daher in der Mensch-Maschine-Kommunikation erhebliche Aufmerksamkeit gewidmet. Jedoch gibt es bislang kaum Lösungsvorschläge dafür, wie die multimodalen Äußerungen eines Systemnutzers – als zeitlich gestreute Perzepte auf getrennten Kanälen registriert – in ihrem zeitlichen Zusammenhang zu rekonstruieren sind. Dieser Beitrag demonstriert erste technische Arbeiten, die rhythmische Muster für die Entwicklung kognitiv motivierter Mittersysteme zwischen Mensch und Maschine ausnutzen.

1 Einleitung

Mit dem verstärkten Eintritt des Menschen in multimediale "virtuelle" Umgebungen finden Formen der nichtverbalen körperlichen Äußerung, insbesondere Gesten, als Mittel der Informationsübermittlung an maschinelle Systeme starkes Interesse. Untersucht werden in jüngerer Zeit auch 'koverbale' Gesten, also Gesten, die sprachliche Äußerungen mehr oder weniger spontan begleiten, z.B. wenn man auf einen Gegenstand zeigt ("dieses Rohr") oder eine Drehrichtung ("so herum") signalisiert. Es ist leicht erkennbar, daß eine derartige Eingabeform für Multimedia-Systeme, wie sie heute schon im virtuellen Entwurf eingesetzt werden, erheblichen Komfortgewinn erbringen könnte.

Als eine Herausforderung stellt sich dabei die *multimodale Integration*, insbesondere die zeitliche Kopplung der beiden komplementären Modalitäten gesprochener Sprache und Gestik: Die von der Natur her multimodalen Äußerungen eines Systemnutzers werden als nebenläufige Sprach- und Gestenperzepte auf getrennten Kanälen technisch registriert und müssen für die Steuerung von Anwendungen zusammengeführt und interpretiert werden. Bei der Vorverarbeitung der in der Signalerfassung aufgenommenen Meßdaten kommt es zu Verzögerungen, und die Zeitkonstanten dieser Prozesse sind verschieden, das heißt, die zentrale Verfügbarkeit von Informationen aus der Signalvorverarbeitung ist zeitlich gestreut. Um für die Interpretation der Meßergebnisse den inhaltlichen Zusammenhang wieder herzustellen, ist also zunächst ihr zeitlicher Zusammenhang zu ermitteln. Technische Verfahren müssen das Zeitverhalten schon deshalb rekonstruieren, damit die Integration des Zeichenhaften (z.B. Zeigegeste) mit dem Signalgehalt (z.B. Zeigevektor im Moment des Zeigens) gelingen kann. Für das Problem der zeitlichen Integration multimodaler Eingaben haben sich bislang keine befriedigenden Lösungen finden lassen.

Anhaltspunkte ergeben sich aber durch Forschungsbefunde aus den Humandisziplinen, die zeigen, daß das menschliche Kommunikationsverhalten durch signifikant 'rhythmische' Muster geprägt ist (Condon, 1986). Wenn eine Person spricht, bewegen sich oft viele Teile des Körpers (Arme, Finger, der Kopf etc.) zur gleichen Zeit und in enger zeitlicher Kopplung (Selbstsynchronität). Die Ausführung einer Geste läßt sich in mehrere Phasen unterteilen, von denen die expressive Phase (Stroke) die wichtigste ist. Der Stroke ist häufig durch einen abrupten Halt gekennzeichnet, der mit den dabei gesprochenen Wörtern zeitlich in enger Beziehung steht. Sprache und Körperbewegungen zeigen dabei charakteristische Periodizitäten; z.B. finden sich in allen germanischen Sprachen – bei flüssigem Sprechen – Korrelationen zwischen (durch zeitliche Dehnung) betonten Silben und einhergehenden Gesten-Stroke. Experimente haben ergeben, daß ein betontes Wort in der Regel nicht vor dem Stroke der koverbalen Geste geäußert wird; der Stroke tritt kurz zuvor oder spätestens

mit dem betonten Wort auf (McNeill, 1992). Auch in der sprachlichen Äußerung allein lassen sich rhythmische Akzentuierungen beobachten, die sich im Timing des Sprechens äußern (Kien & Kemp, 1994); (Fant & Kruckenberg, 1996). Ebenfalls ist verschiedentlich beobachtet worden (Condon, 1986); (McClave, 1994), daß die Äußerungsrhythmik eines Sprechers vom Hörer in körperlichen Reaktionen übernommen wird (Interaktionssynchronität).

Ähnlich wie die rhythmische Koordination der Gliedmaßen bei der Lokomotion (Schöner & Kelso, 1988) werden kommunikativen Rhythmen als koordinative Strategie des menschlichen Äußerungs- und Wahrnehmungsapparats gedeutet. Rhythmen scheinen eine Art "Taktschläge" bereitzustellen, die die Synchronisation von Körperbewegung und gesprochener Sprache bewerkstelligen. Durch erwartbare Periodizitäten bringen sie quasi vereinzelbare Prozeßeinheiten hervor, die dem Rezipienten das Segmentieren des übertragenen Signals erleichtern (Martin, 1979). Weitere Hinweise geben Untersuchungen zu temporalen Kontrollmechanismen für die Wahrnehmung und Bewußtseinsbildung im menschlichen Gehirn (Pöppel, 1997). Danach werden aufeinanderfolgende Wahrnehmungszustände zu größeren, bis zu drei Sekunden langen, Einheiten gebündelt. Dies gilt insbesondere für Verbindungen zwischen den verschiedenen Sinnesmodalitäten. Ebenfalls gibt es Hinweise, daß die Abfolge von absichtlichen Bewegungen, zu denen auch die meisten Gesten zählen, zeitlich strukturiert und bis zu einem Zeitraum von 2 bis 3 Sekunden vorausgeplant ist.

2 Ansätze für die multimodale Integration mit Rhythmen

In der Arbeitsgruppe Wissensbasierte Systeme der Universität Bielefeld werden seit mehreren Jahren Möglichkeiten der Gestenerkennung für Mensch-Maschine-Schnittstellen und der multimodalen Integration von Gestik und Sprache erforscht (z.B. Wachsmuth, 1999a). Die im Einleitungsabschnitt geschilderten Beobachtungen führten uns zu dem Gedanken, die Analyse kommunikativer Rhythmen zur Verbesserung der Leistungsfähigkeit technischer Mittlersysteme zwischen Mensch und Maschine auszunutzen. Gesprochene Sprache und Gestik sind zunächst einmal essentiell kontinuierliche Prozesse. Vor einer semantischen Analyse übermittelter Information sind also die folgenden logistischen Probleme zu lösen (Srihari, 1995):

(1) *Das Segmentierungsproblem:* Wie sind die Prozeßeinheiten zu determinieren, die das System in einem Zyklus verarbeiten soll?

(2) *Das Korrespondenzproblem:* Wie sind die Querbezüge zwischen den Modalitäten Gestik und Sprache zu determinieren?

Unter der Annahme, daß ein grundlegender Takt im Äußerungsverhalten des Menschen besteht, könnte durch Verwertung von Segmentierungshinweisen, wie Gesten-Stroke und Sprechtakt, der kommunikative Rhythmus systemseitig reproduziert und u.U. antizipiert werden. Dies würde dabei helfen, die Korrespondenzen der zeitlich gestreuten Sprach- und Gestenperzepte wieder herzustellen und dadurch die semantische Analyse multimodaler Information erleichtern.

In einem ersten technischen Ansatz im Projekt VIENA ("Virtuelle Entwurfsumgebung und Agenten") wurde das Prinzip der kommunikativen Rhythmen zur Bestimmung von zusammengehörigen Worten und Zeigegesten ausgenutzt (Lenzmann, 1998); siehe auch (Wachsmuth, 1999b). Das in VIENA realisierte System kann z.B. bei Instruktionen wie "put - <Geste> this - computer - on - <Geste> that - desk" die über Spracherkenner und Datenhandschuh registrierten Eingaben zusammenführen, um entsprechende Änderungen in einer computergrafisch visualisierten Szene zu berechnen. Die Korrespondenz von Zeigegesten und bei der Sprachanalyse determinierten 'Gestenplätzen' (das sind Informationsplatzhalter, die Erwartungen bezüglich ergänzender Objekt- und Richtungsspezifikationen formalisieren) ist u.a. durch die zeitliche Nähe geleitet: Je kleiner der Abstand, desto besser passen beide Teile zusammen. Der Integrationsinstanz des Systems, dem sog. Koordinator, ist ein 2-Sekunden-Rhythmus aufgeprägt; er wird durch das erste Ansprechen des Mikrophons angestoßen und sorgt dafür, daß die in einem solchen "Takt" registrierten Ereignisse a priori als zusammengehörig betrachtet werden. Der angestoßene Rhythmus ("swing") klingt aus ("subside"), wenn keine Eingabeereignisse mehr registriert werden, und geht bis

zu einer Folgeinstruktion in einen Wartezustand ("wait") über. Der Ansatz unterstützt zudem offene Eingaben, indem nach Ablauf eines Takts automatisch eine Integration der Ereignisse vorgenommen wird und so keine explizite Markierung des Eingabeendes erforderlich ist; die Segmentierung erfolgt allein durch den im Koordinator evozierten Rhythmus.

Die Erstellung eines umfassenden Systemprototyps, der von der Erkennung komplexerer Gesten über die Sprach-Gestik-Integration bis zur Anbindung an eine Zielapplikation des virtuellen Konstruierens reicht, ist das Kernziel des SGIM-Projekts ("Sprach- und Gesten-Interfaces für Multimedia"). In den hier entwickelten verfeinerten Ansätzen ist die Methode `rhythmInfo` für die Verarbeitung von rhythmisch strukturierter Information zuständig. Sie erwartet als Parameter den Zeitpunkt eines Taktschlags, der die Grenze eines semantischen Segments andeutet. Alle Signalperzepte, deren Assertionszeit älter als die Taktzeit ist, werden aus dem Arbeitsgedächtnis des Systems entfernt; es wird damit zyklisch von nicht mehr relevanter Information befreit. Rhythmusinformation wird mit der Message-Klasse `RhythmMessage` im System kommuniziert. Als einzige Komponente enthält die Klasse einen Zeitstempel, der den genauen Zeitpunkt des Taktschlags mitteilt; er liegt je nach Signaltyp innerhalb oder direkt am Anfang einer neuen semantischen Einheit. Als Reaktion auf eine `RhythmMessage` wird im derzeitigen System die Methode `rhythmInfo` des Integrators aufgerufen. `RhythmMessage` ist die Oberklasse aller Nachrichtenklassen, die Rhythmusinformation versenden. Die Teilklasse `GestureSegmentationCue` ist für Gestensegmentierungshinweise zuständig. Als Komponente enthält sie eine Markierung, die angibt, ob der Zeitstempel die Grenze zweier semantischer Einheiten oder die expressive Phase einer Geste bezeichnet. Als Beispiel eines Gestensegmentierungshinweises arbeiten wir zum Beispiel mit einer Klasse `HandTensionMessage`, die sich zunutze macht, daß zwischen je zwei ausgeprägten Gesten sich die Hand kurzfristig entspannt, was über die Meßsignale eines Datenhandschuhs feststellbar ist. Ein regelbasiertes Rahmensystem, in dem die zeitliche Integration symbolischer Information aus unterschiedlichen Modalitäten realisiert wird, wird in (Sowa, Fröhlich & Latoschik, 1999) beschrieben.

3 "Rhythm is the key"

Sieht man sich die anfangs angerissenen, durchaus frappierenden Befunde zu den hier thematisierten "kommunikativen Rhythmen" näher an, zeichnet sich das Bild ab, daß in der Kommunikation zeitlich-strukturellen – und damit auch rhythmischen – Merkmalen ein ebenso großer Stellenwert wie der semantischen Informationsverarbeitung einzuräumen ist. In der Mensch-Maschine-Kommunikation haben solche Erkenntnisse bislang kaum Eingang gefunden. Unsere Untersuchungen lassen hoffen, daß Rhythmus der Schlüssel für einige schwierige Probleme, insbesondere in der multimodalen Kommunikation sein könnte. Auf jeden Fall eröffnet sich hiermit ein Diskursthema, das spannende Forschungsfragen für die kognitiven Disziplinen verspricht.

Literatur

- Condon, W.S. (1986). Communication: Rhythm and Structure. In J. Evans and M. Clynes (Eds.): *Rhythm in Psychological, Linguistic and Musical Processes* (pp. 55-77). Springfield, Ill.: Thomas.
- Fant, G. & Kruckenberg, A. (1996). On the Quantal Nature of Speech Timing. *Proc. ICSLP-96*, pp. 2044-2047.
- Kien, J. & Kemp, A. (1994). Is speech temporally segmented? Comparison with temporal segmentation in behavior. *Brain and Language* 46: 662-682.
- Lenzmann, B. (1998). *Benutzeradaptive und multimodale Interface-Agenten*. Dissertationen der Künstlichen Intelligenz, Bd. 184, Sankt Augustin: Infix.
- Martin, J.G. (1979). Rhythmic and segmental perception. *J. Acoust. Soc. Am.* 65(5): 1286-1297.

- McClave, E. (1994). Gestural Beats: The Rhythm Hypothesis. *Journal of Psycholinguistic Research* 23(1), 45-66.
- McNeill, D. (1992). *Hand and Mind: What Gestures Reveal About Thought*. Chicago: University of Chicago Press.
- Pöppel, E. (1997). A hierarchical model of temporal perception. *Trends in Cognitive Science* 1(2), 56-61.
- Schöner, G. & Kelso, J.A.S. (1988). Dynamic pattern generation in behavioral and neural systems. *Science*, 239: 1513-1520.
- Sowa, T., Fröhlich, M. & Latoschik, M. (1999). Temporal symbolic integration applied to a multimodal system using gestures and speech, presented at GW'99: 3rd Internat. Gesture Workshop, 17-19 March, 1999, Gif-sur-Yvette.
- Srihari, R.K. (1995). Computational models for integrating linguistic and visual information: a survey. *Artificial Intelligence Review* 8: 349-369.
- Wachsmuth, I. (1999a). Mensch-Maschine-Kommunikation mit Gestik und Sprache, ersch. in: W.-D. Miethling & J. Perl (Hrsg.), *Sport und Informatik VI* (S. 167-178). Köln: Sport und Buch Strauss.
- Wachsmuth, I. (1999b). Communicative rhythm in gesture and speech, presented at GW'99: 3rd Internat. Gesture Workshop, 17-19 March, 1999, Gif-sur-Yvette, to appear (Springer).