

Published in final edited form as:

Intell Virtual Agents. 2014 January 1; 8637: 438-448. doi:10.1007/978-3-319-09767-1_55.

A Virtual Therapist for Speech and Language Therapy

Sarel van Vuuren^{1,2} and Leora R. Cherney^{3,4}

Sarel van Vuuren: sarel@colorado.edu; Leora R. Cherney: lcherney@ric.org

¹University of Colorado Boulder, Institute of Cognitive Science, Boulder, CO

²University of Colorado Anschutz Medical Campus, School of Medicine, Department of Physical Medicine and Rehabilitation, Denver, CO

³Rehabilitation Institute of Chicago, Center for Aphasia Research and Treatment, Chicago, IL

⁴Northwestern University, Feinberg School of Medicine, Department of Physical Medicine and Rehabilitation, Chicago, IL

Abstract

A virtual therapist (VT) capable of modeling visible speech and directing speech and language therapy is presented. Three perspectives of practical and clinical use are described. The first is a description of treatment and typical roles that the VT performs in directing participation, practice and performance. The second is a description of techniques for modeling visible speech and implementing tele-rehabilitation. The third is an analysis of performance of a system (*AphasiaRx*TM) for delivering speech and language therapy to people with aphasia, with results presented from a randomized controlled cross-over study in which the VT provided two levels of cuing. Compared to low cue treatment, high cue treatment resulted in 2.3 times faster learning. The paper concludes with a discussion of the benefits of speech and language therapy delivered by the VT.

Keywords

Virtual agents; virtual therapists; speech and language therapy; aphasia treatment; tele-rehabilitation; tele-care; speech animation

1 Introduction

For persons with aphasia (PWA), virtualization of treatment is critical given the prevalence of the disorder, effect it has on their ability to communicate, and general lack of access to long-term treatment options. This paper describes the role, implementation and performance of a virtual therapist (VT) for delivering speech and language therapy.

Aphasia is an acquired multi-modality disturbance of language, resulting from focal damage to the portions of the brain responsible for language. It impairs, to varying degrees, understanding and expression of oral language, reading and writing. It occurs with a wide range of neurologic disorders including stroke, brain tumor, cerebral trauma, intracranial surgical procedures, and degenerative neurological disorders [1]. An estimated 80,000 new cases of aphasia per year in the United States result from stroke alone, while the prevalence

in the country overall has been estimated to exceed one million people [2]. Effects are chronic and potentially long-term, with reduced language skills limiting participation in social, vocational, and recreational activities.

Services for persons with aphasia (PWAs) have been drastically shortened. Legislation has seriously curtailed the amount of treatment a PWA may receive after hospitalization. Often they may be eligible for only a limited number of treatment sessions for a limited period of time. In some cases, they may not receive any treatment for their communication disorder following their acute hospitalization [3]. Reduced resources (e.g. transportation difficulties, therapist shortages in rural areas) further limit available services. Given this state of affairs, it is imperative that clinicians provide treatment that results in the greatest improvement in the shortest amount of time— motiving a shift within the field to tele-rehabilitation for anywhere-access, and as described in this paper, reliance on a VT to maintain guidance, encouragement and fidelity.

1.1 Treatment with a Virtual Speech Therapist

Script treatment, whether delivered by human [4, 5], or computer [6, 7, 8], has been shown to be highly effective for treating PWAs. In script treatment, therapy is based on the oral production of scripts, which are short functional dialogs structured around communication of everyday activities. A computer program called *AphasiaScripts*TM that we developed in 2004 is available from www.ric.org/aphasia, but its use of a proprietary VT with visible speech based on motion capture variable length concatenative synthesis [9], with resulting large footprint and restriction to computer only use, meant it was not suited for telerehabilitation across devices. Recently, we developed a new system called *AphasiaRx*TM, designed to be more portable and functional, with a new VT that works across devices, and built-in web monitoring, scheduling and communication technologies. Fig. 1 shows a screen image.

The new system is designed to be easy to use by people with developmental and acquired disabilities, who might have fine motor skills difficulties. Navigation is done with a space bar, and in the case of a tablet, with the whole screen as a button. Activities are sequenced, leveled, and adaptive where possible, with guidance and feedback provided interactively by the VT. Because the treatment is repetitive, VT responses are limited to short cues to maximize practice time and treatment intensity.

Treatment is structured around three phases of script practice:

- 1. The PWA listens to a script while it appears on the screen. The VT is visible and using a high quality recorded female voice, reads aloud, articulating carefully. Each word is read and highlighted on the screen.
- 2. Each sentence or conversation turn is read aloud twice in unison with the VT. Problem words can then be practiced repeatedly by clicking on them. Finally, the PWA independently reads aloud the whole sentence while it is recorded by the system, and is given the option to listen to the recording or practice the sentence again.

3. The conversation is practiced with assistance provided by the VT. Maximum support includes seeing written words, hearing the VT's voice during choral speaking, and watching her oral-motor movements. As the PWA masters the script, he or she can remove cues one by one, so that eventually practice simulates real conversation.

For the speech-language pathologist (SLP), built-in tools make authoring new scripts and treatment plans, assigning them to their clients, and monitoring their performance easy. Scripts can be typed into the program and recorded. Each spoken word is aligned automatically with each written word. Pause time between sentences and the amount of time allotted for each utterance can be adjusted to provide individualized optimal speaking time for each. Different scripts, treatment plans, and schedules can be assigned to individual clients with a built-in calendar. Data are automatically captured, summarized and available for analyses. During practice, key strokes, item responses and response times are recorded and summarized to provide daily and weekly logs showing how time was spent in each part of the program. Recordings of the speech attempts made by the PWA provide the SLP with additional means of assessing change over time.

1.2 Related Approaches in Other Areas

Many studies support the usefulness of intelligent virtual agents (IVAs) that are conversational, embodied, and/or animated for activities such as interaction, motivation and tutoring. E.g., see [10, 11, and 12] for an overview. Indeed, IVAs seem to add a "human dimension" to the computer interface that improve, direct and encourage participation, communication and accessibility. In our own research, IVAs have helped users in educational applications ranging from reading instruction [13] to science tutoring, and VTs have helped users remediate and rehabilitate speech-language difficulties in Parkinson's disease [14] and aphasia [7, 15, 16, 17]. For aphasia in particular, while video based treatment methods exist, e.g. [18], our approach to script treatment described here is unique in its use of high-fidelity IVAs.

2 Modeling and Implementation

There are many existing approaches for animating agents and speech to varying degrees of accuracy, with tradeoffs in computing, cost, flexibility and availability. Our goal is to implement therapy across a range of devices, yet maintain reasonable accuracy, simplicity, consistency, and control. This section explains briefly how this is done.

2.1 Visual Animation

For local or remote control of the VT's visible speech on computer, web and mobile devices, an approach that decouples modeling and animation is needed. To do this, we use a code book to capture visual variation. We assume that the visual realization of individual sounds vary depending on the phonetic context and manner of speaking, and that the temporal evolution depends predominantly on physiological and acoustic constraints governing movement of the oral-motor articulators during acoustic realization— though at different time scales and rates [19]. Using a code-book ensures that the link between speech and animation is in the choice of timing and indices.

Animation starts with a phoneme sequence and corresponding time-aligned audio generated using a TTS engine or pre-recorded speech aligned with a speech recognizer [13]. The phoneme sequence is then mapped using a code-book to 3D morphs (or 2D images, depending on the application), so that individual phonemes correspond to a sequence of code-book indices in a one-to-many mapping, where the mapping depends on the dynamic realization of the phoneme as well as co-articulation effects.

Code-book entries represent the extremes of visual realizations of parts of phonemes with indices spanning a low 5-dimensional 'visual articulation' space, shown conceptually in Fig. 2. When rendered, the entries look somewhat like the images in Fig. 3.

The space was designed manually using phonological, linguistic and morphological considerations to cover the range of speech locomotion observed in prototypical image sequences and videos that we collected across selected permutations of normally spoken phonemes, syllables and words. For this purpose, we assumed that the dynamic evolution of visible speech is relatively continuous and low-dimensional, so that it can be modeled as a process that traces a continuous trajectory within the space.

Rather than trying to model the articulation space parametrically, the code-book quantizes a convex region of the space on a non-linear sparse grid, allowing the trajectory contained within to be approximated as a sequence of code-book indices. Convexity, monotonicity and smoothness of the space are part of the design so that interpolation and extrapolation can be used. Because indices correspond to grid points, co-articulation, speech rate, and speech effort are modelled by applying transformations directly to the sequences. The transformations are applied at run time and are relatively simple smoothing, blending, scaling and translational functions.

Finally, the resultant sequence of indices are used to recover from the code-book a sequence of high-dimensional morphs (or static images, depending on the application) of the mouth, lips, teeth and tongue, to be interpolated and rendered at the device supported frame rate (Fig. 3). For computer or mobile devices, OpenGL is used to render 3D models, whereas for web applications, HTML5 and WebGL or pure images are used to render 3D models or 2.5D approximations, respectively. We chose this implementation over existing methods (e.g., www.facefx.com), because we wanted full control of the implementation across devices in order to optimize it for tele-rehabilitation.

The approach is related to previous work, but also differs in some aspects. Facial expressions, gaze, and idle animations are modelled similar to [9], though not generated from speech as in [20]. The articulation space is a composite blend of various articulation parameters, making it somewhat different from other parametric approaches [21, 22], but similar in that it allows explicit manipulation of the visible speech. The code-book is designed to allow physiological and articulatory phonetic manipulation of the visible speech, making it different from other data-driven approaches such as video/image stitching [12] and unit-concatenation [9], but similar in that it attempts to model variation using code-book entries, though its entries differ from visemes. The coding scheme, while suitable for

transmission, captures visual information at a more linguistic level than the facial animation parameters in the M-PEG 4 standard.

The visible speech produced by the agent has been qualitatively reviewed by speech language pathologists, linguists and audiologists, and we have controlled experiments underway to test the accuracy clinically. Thus, we cannot currently provide specific results regarding visible accuracy or speech intelligibility. Instead, in Section 3 we demonstrate ecological validity with a randomized controlled cross-over study for a speech and language therapy treatment that compared two conditions, one in which the VT modeled speech production, and one with little modeling of speech production.

2.2 Tele-Rehabilitation

Implementation follows similar considerations for fidelity and access as [13, 15, 16], but is more flexible and sophisticated in that the system described here can work across a number of different devices. It can operate either as a standalone native application on a computer, tablet, or mobile device, or as a web application in a cloud-based client-server configuration suitable for tele-rehabilitation (Table 1). For the latter, communication, synchronization and access control between the clients (applications) and server are done entirely with the HTTP/S protocol and asynchronous calls.

While this paper focuses on an application for aphasia treatment, the VT and system is suitable for applications that extend beyond virtual speech and language therapy. For example, the architecture's underlying capabilities for web and mobile rehabilitation have been used between Boulder, Denver, and Toronto in a near real-time non-linear context aware prompting system (NCAPS) for people with cognitive disabilities, helping them to navigate and overcome procedural challenges in the workplace [23].

3 Clinical Application

A randomized controlled cross-over study using the VT was conducted to investigate the effect of high or low cuing on treatment outcomes over time. It was approved by the IRB of Northwestern University and done at the Rehabilitation Institute of Chicago.

3.1 Methodology

Eight participants were recruited and randomized to receive intensive computer-based script training differing in the amount of high or low cuing provided during treatment. Participants had chronic aphasia resulting from a single left hemisphere stroke (age range 25–66 years, m=52.0, sd = 14.0; time post-onset range 8–59 months, m=26.4, sd=19.2), with severities mild-moderate to severe as measured by the Aphasia Quotient of the Western Aphasia Battery-Revised [24] (range 28.1–80.1, m=58.0, sd=18.5).

In the high cue treatment condition, participants could hear the VT during listening, choral reading and reading aloud, with auditory cues (therapist speaking) and visual cues (therapist's mouth movements) available at the start, during and after practice. In the low cue condition, they received visual and auditory cues when listening to the script being read

aloud *initially*, and *after* practice; but did not receive auditory and visual support *during* sentence practice, working instead from the written cues only.

Confounds were controlled for with a cross-over design with three weeks of treatment in one condition, followed by a three week washout period and three weeks of treatment in the second condition; random assignment, and stratification of participants to treatment conditions; control for treatment time; and design and use of similar though different scripts in each condition with scripts matched for length, morphological, phonological and grammatical complexity. Scripts consisted of 10 sentences each.

Scripts were practiced 6 days per week, for up to 90 minutes per day in 3 sessions. Practice and testing were done at home on a loaned laptop with sessions scheduled automatically using the system calendar. Participant responses were measured over time: in 3 sessions before treatment started to establish a baseline, then in separate sessions twice per week during each of the three treatment weeks to assess learning, and finally, once in the week after treatment finished to assess treatment effectiveness.

Performance was measured by averaging the sentence level word accuracy of participants' production of 10 sentences during each assessment session. Sentences were 10 words in length. Accuracy of words were rated using a previously validated 6 point scale and the overall session score expressed on a scale from 0 to 100%. For the time period and measures reported here, a total of 16,000 words (8 participants \times 2 conditions \times 10 sessions \times 10 sentences \times 10 words) were collected, recorded and scored. Interrater reliability on a 10% subset of sentences was 0.94 (Pearson's r).

In previous work we discussed the implications of this study for speech-language pathology and reported pre- and post-treatment results conditioned on aphasia severity [17]. There, we used Cohen's d to compute effect sizes on gains, with d>0.8 considered a large effect. When combining the high cue and low cue conditions (16 samples, i.e. including all participants for N=8 and two treatment conditions), the effect size for computer treatment was large and statistically significant (d=1.5, p<0.0125).

In the following, we report new results showing the differential effects of high cue and low cue treatment *over time* when delivered by the VT described in Section 2.

3.2 Effect of Cuing on Treatment Over Time

Figs. 4a and 4b show the effect of treatment over time. Gain in scores (% accuracy) over baseline are shown before, during and after treatment. Gains in both high and low cue treatment conditions increased rapidly once treatment started.

Fig. 4a shows the effect of treatment and cumulative practice hours for one participant. Practice time in both conditions were matched, and congruent with gains. At the end of treatment gains leveled off. For a perfect score, the participant's maximum attainable gains were limited to 36.6% in the high cue condition and 38.3% in the low cue condition, suggesting possible ceiling effects during the latter half of the treatment period. Limited participant pre-training before the treatment period can also be seen.

Fig. 4b shows average gain scores (% accuracy) over baseline for all 8 participants, with the effects of high cue and low cue treatment shown separately over time. Values were computed by averaging samples for participants in each condition over both the time and value axes. The maximum gain that could be attained depended on the participant's starting level and ranged from 14.8% to 38.3% in the low cue treatment condition, and 15.2% to 37.3% in the high cue treatment condition.

Of the eight participants, six learned faster in the high cue condition. Indeed, gains supported findings from a post-treatment survey reported in [17] where 6 of 8 participants said they liked high cue treatment better.

High cue treatment on average led to faster learning and higher gains, with an estimated overall speedup of L/H=2.3 using the averages from Fig. 4b. This meant that when participants received high cue treatment, they on average learned more than twice as fast as when they received low cue treatment, reaching the same level of proficiency in less than half the time (1.3 weeks compared to 2.9 weeks from start of treatment).

Preliminary analyses with mixed effects models – often used in pharmacological studies of treatment and dosage response over time [25], showed the effect is statistically significant. A logistic mixed effects model fitted to the gains for all 8 PWAs and conditions, showed that cue level differentially and significantly affected gain over the treatment period ($\chi^2(1)$ =10.98,p<0.0009) [paper in preparation]. From the model, the predicted speedup for the difference in time to reach the same level of gain between the two conditions was L/H=2.34, similar to the non-parametric estimate from Fig. 4b.

4 Discussion and Conclusions

The paper described a virtual therapist (VT) for delivering speech and language therapy to persons with aphasia (PWA). Three different perspectives focusing on role, implementation and performance were provided, namely: how the VT directed participation, practice and performance in much the same way an SLP does; modeling and implementation considerations for visible speech and tele-rehabilitation; and evidence of ecological validity, with results presented from a randomized controlled cross-over study in which the VT provided two levels of cuing, with significantly faster learning (2.3×) for high cue treatment compared to low cue treatment.

Caveats to the results were that the sample size was small and caution should be exercised with interpretation. The study did not attempt to separate the effects of auditory and visual cuing—work that might be explored in future research. While the study did not compare treatment delivered by the VT to treatment delivered by an expert SLP, the overall effect size for the computerized treatment was large (d=1.5) and similar to what would be expected when treatment is delivered by an SLP. Nevertheless, the results reported here were for acquisition only without consideration of possible interactions with maintenance and generalization post-treatment.

In summary, the paper showed that for persons with aphasia, receiving treatment in an ecologically valid real-world setting delivered by a VT that provides more cues than not, can lead to faster learning.

Acknowledgements

Supported by the National Institute on Deafness and Other Communication Disorders, National Institutes of Health, Award 1R01DC011754 (to S.V.V. and L.R.C.). Animation research supported by University of Colorado Boulder, and platform research supported in part by the National Institute on Disability and Rehabilitation Research, U.S. Department of Education, Award H133E090003 (to C. Bodine, www.rerc-act.org). The content is solely the responsibility of the authors and does not necessarily represent the official views of the funding organizations. Endorsement by the Federal Government should not be assumed. The authors do not have a financial interest in the product. Thanks are extended to Julia Carpenter, Rachel Hitch, Rosalind Hurwitz, Rosalind Kaye, Jaime Lee, Anita Halper, Audrey Holland, Michael C. Mozer, Nattawut Ngampatipatpong, Robert Bowen and Taylor Struemph.

References

- Albert, ML.; Goodglass, H.; Helm, NA.; Rubens, AB.; Alexander, MP. Clinical Aspects of Dysphagia. New York: Springer-Verlag Wien; 1981.
- 2. NIH Pub No. 97-4257: National Institute on Deafness and Other Communication Disorders, Facts Sheet: Aphasia. Bethesda, MD: Author; 2008.
- 3. Elman, RJ.; Simmins-Mackie, N.; Kagan, A. Clinical services for aphasia: Feast or famine?. Paper presented at the Annual Meeting of the American Speech-Language-Hearing Association; Chicago, Illinois. 2003.
- Holland, A.; Milman, L.; Munoz, M.; Bays, G. Scripts in the management of aphasia. Paper presented at the World Federation of Neurology, Aphasia and Cognitive Disorders Section Meeting; Villefranche, France. 2002.
- 5. Youmans GL, Holland AL, Munoz M, Bourgeois M. Script training and automaticity in two individuals with aphasia. Aphasiology. 2005; 19:435–450.
- Cherney LR, Halper AS. Novel Technology for Treating Individuals with Aphasia and Concomitant Cognitive Deficits. Topics in Stroke Rehabilitation. 2008; 15(6):542–554. [PubMed: 19158062]
- 7. Lee JB, Kaye RC, Cherney LR. Conversational script performance in adults with non-fluent aphasia: Treatment intensity and aphasia severity. Aphasiology. 2009; 23(7):885–897.
- 8. Manheim LM, Halper AS, Cherney L. Patient-Reported Changes in Communication after Computer-Based Script Training for Aphasia. Archives of Physical Medicine and Rehabilitation. 2009 Apr; 90(4):623–627. [PubMed: 19345778]
- 9. Ma J, Cole R, Pellom B, Ward W, Wise B. Accurate visible speech synthesis based on concatenating variable length motion capture data. IEEE Transactions on Visualization and Computer Graphics. 2006; 12(2):266–276. [PubMed: 16509385]
- 10. Moreno R, Mayer R, Spires H, Lester J. The Case for Social Agency in Computer-Based Teaching: Do Students Learn More Deeply When They Interact with Animated Pedagogical Agents? Cognition and Instruction. 2001; 192(2):177–213.
- 11. Gratch J, Rickel J, André E, Badler N, Cassell J, Petajan E. Creating interactive virtual humans: Some assembly required. IEEE Intelligent Systems. 2002; 17(4):54–63.
- 12. Cosatto E, Ostermann J, Graf HP, Schroeter J. Lifelike talking faces for interactive services. Proceedings of the IEEE: Special Issue on Human-Computer Multimodal Interface. 2003 Sept. 91(9):1406–1429.
- Van Vuuren S. Technologies that power pedagogical agents. Educational Technology. 2007; 24(1):
 4–10
- 14. Cole R, Halpern A, Ramig L, Van Vuuren S, Ngampatipatpong N, Yan J. A Virtual Speech Therapist for Individuals with Parkinson Disease. Educational Technology. 2007; 24(1):51–55.
- 15. Cherney, L.; Babbit, E.; Kwang-Youn, K.; Van Vuuren, S.; Ngampatipatpong, N. Aphasia Treatment over the Internet: A Randomized Placebo-Controlled Clinical Trial; Clinical Aphasiology Conference, Ft.; Lauderdale, FL. May 31–June 4; 2011.

 Cherney L, Van Vuuren S. Telerehabilitation, Virtual Therapists and Acquired Neurologic Speech and Language Disorders. Seminars in Speech and Language. 2012 Aug; 33(3):243–57. PubMed PMID: 22851346. [PubMed: 22851346]

- Cherney L, Kaye R, Van Vuuren S. Acquisition and Maintenance of Scripts in Aphasia: A Comparison of Two Cuing Conditions. American Journal of Speech-Language Pathology. 2014; 23:S343–S360. [PubMed: 24686911]
- 18. VAST speech aid. SpeakinMotion. 2014 http://www.speakinmotion.com.
- Yang HH, Van Vuuren S, Sharma S, Hermansky H. Relevance of time-frequency features for phonetic and speaker-channel classification. Speech Communication. 2000 May; 31(1):35–50.
- Marsella S, Xu Y, Lhommet M, Feng A, Scherer S, Shapiro A. Virtual character performance from speech. Proceedings of the 12th ACM SIGGRAPH/Eurographics Symposium on Computer Animation. 2013:25–35. ACM.
- 21. Cohen, M.; Massaro, D. Modeling coarticulation in synthetic visual speech. In: Thalmann, N.; D, T., editors. Models and Techniques in Computer Animation. Springer-Verlag; 1994. p. 141-155.
- Pelachaud, C. Visual Text-to-Speech, in MPEG4 Facial Animation The standard, implementations and applications. Pandzic, Igor S.; Forchheimer, Robert, editors. John Wiley & Sons; 2002.
- 23. Melonis M, Mihailidis A, Keyfitz R, Grzes M, Hoey J, Bodine C. Empowering Adults With a Cognitive Disability Through Inclusion Of Non-Linear Context Aware Prompting Technology (N-CAPS). RESNA Conference Proceedings. 2012 Jun 30.
- 24. Kertesz, A. Western Aphasia Battery-Revised. San Antonio, TX: PsychCorp.; 2007.
- 25. Pinheiro, JC.; Bates, DM. Mixed-Effects Models in S and S-PLUS. Springer; 2000.

PROBE SCRIPT

Subway: Can I help you?

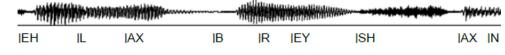
You: I'd like a turkey sandwich.



press spacebar to continue

Fig. 1. Screen image from $AphasiaRx^{TM}$ system with sentence for VT and PWA turn displayed.

Input: Time-aligned phoneme sequence



Encoding and decoding step:

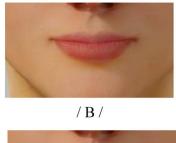
Region in articulation space

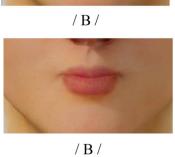
Trajectory

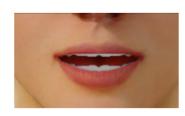
Output: Code-book entries to be rendered as morphs or images

Fig. 2.Conceptual depiction of speech trajectory evolving in visual articulation space and quantization to code-book spanning the space. Actual space is 5 dimensional. The code-book maps a time-aligned phoneme sequence to a corresponding sequence of morphs or images to be rendered.









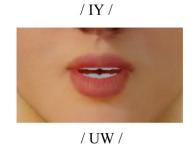


Fig. 3. Screen image of the animated agent and close up of visible speech showing co-articulation for the consonant-vowel pairs $\/B$ IY $\/$ and $\/B$ UW $\/$ in the words 'beet' and 'boot', respectively.

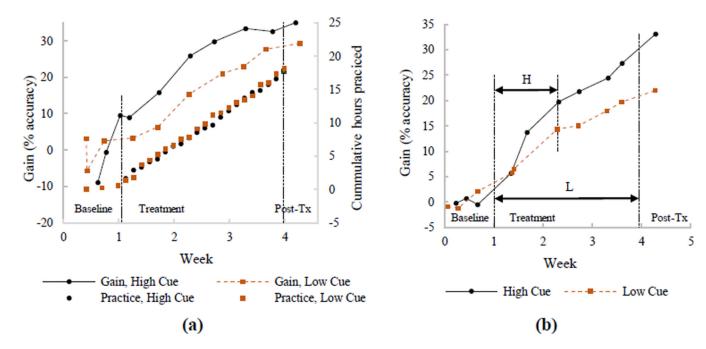


Fig. 4.

Gain (in % accuracy) over baseline for high and low cue treatment conditions. Gains increased rapidly once treatment started, with gains for high cue increasing faster than low cue treatment. (a) Gains and cumulative hours practiced for one participant. (b) Average gains for N=8 participants. Participants receiving high cue compared to low cue treatment reached the same level of gain faster, in about half the time.

Table 1

The system client-server architecture suitable for tele-rehabilitation.

Client-side User	Asynchronous HTTP/S	Server-side Control
Interface (Apps)	Communication	and Data (Cloud)
Computer Web Tablet/Mobile	User and sensor inputs \rightarrow \leftarrow System outputs: Animation control, audio, video, HTML	Task/state management, monitoring, configuration, application and access control, application and user data