

A dialogue-based software architecture for gamified discrimination tests

Antonio Origlia
Dept. of Information
Engineering
University of Padua
antonio.origlia@dei.unipd.it

Piero Cosi
Institute of Cognitive Sciences
and Technology (CNR-ISTC)
piero.cosi@pd.istc.cnr.it

Antonio Rodà
Dept. of Information
Engineering
University of Padua
roda@dei.unipd.it

Claudio Zmarich
Institute of Cognitive Sciences
and Technology (CNR-ISTC)
claudio.zmarich@cnr.it

ABSTRACT

In this work we describe the current stage of development of a software architecture designed to present discrimination tests to pre-school children in the form of gamified tasks. We interpret the problem of administering these tests as a dialogue model using probabilistic rules to generate customised tests on the basis of the child's performance. In the proposed architecture, the dialogue system controls a gaming setup composed of a virtual agent and a robotic companion that needs to be taught how to talk. This learning-by-teaching approach is used to camouflage a phonemes discrimination test that has the added value of being generated at runtime on the basis of the child's performance. We will describe the architectural components involved and we will describe how the dialogue system can make use of linguistic knowledge to generate the discrimination test and administer it by controlling the agents involved in the game.

Author Keywords

Gamification; software architecture; discrimination tests

INTRODUCTION

Phonetic perception abilities are in place and active already in the fetus, and their integrity is necessary for a normal functioning future speech development [12, 23]. Since the ability to discriminate linguistic sounds is associated to the correct acquisition and production of the same sounds, an alteration of the same ability could contribute to the onset of speech and language disorders [2]. For this reason the evaluation of the phonetic discrimination ability is important in order to individuate at-risk subjects, allowing clinicians and caregivers to operate in focused and specific ways. For preschool children (from 3 years-old onward), the paradigms of identification and discrimination are the same as used by adults

[16]. Among several types of discrimination tests, we choose the standard AX or "same-different" procedure. Traditionally, AX tests to evaluate the phonemes discrimination capability of young children are designed as scripts and software traditionally used to administer this kind of test also follows scripts (e.g. [1]). These contain a series of (non-) word pairs presenting phoneme oppositions (i.e. 'pepi / 'pemi) in different syllabic structures (i.e. CV-CV is a disyllabic structure where each syllable has a single heading consonant). The child is given the task to indicate, after listening to the experimenter reading the stimuli, whether the two (non-)words are the same or if they are different. These tests are designed in such a way that consonants presenting a single distinctive trait are opposed at each time (e.g. voiced/unvoiced, sonorant/non-sonorant). Control stimuli are present in such tests as pairs composed by the same word repeated twice and by pairs composed by completely different words. This approach is necessary as it is impossible for a human expert to dynamically select word pairs that comply to a set of very strict constraints. Specifically, each word pair must:

- present opposed consonants that differ in exactly one trait
- syllabic structure must be the same in the two (non-)words
- present the opposition in a precise position in the syllabic structure (e.g. the head consonant of the second syllable)
- the accent must be in the same place in the two (non-)words

Given the young age of the considered subjects, it is necessary to mask the test in a game-like scenario to make it less imposing. Healthy contact with language, in the first years of life, consists of a playful activity where parents and infants engage *protoconversations* made of rhythmical and musical content. This manifests the emotional regulation of *primary inter-subjectivity* [19], where interaction with the caregiver, either reciprocally directed or mediating access to objects of interest for the infant, manifests the typical playfulness often observed in mammals. At 9 months, *secondary inter-subjectivity* arises [22] and the baby's interest moves onto sharing the ways companions use objects as she starts to interact with the material world in a more informed way. The caregivers' language also shifts, in this phase, from questions

and rhetorical comments to instructions and informative comments to support the baby's interest in participating to a task [10]. This is "[...] the start of cultural information transfer between generations" [20, p. 74]. Playful behaviour adapts to new roles as the child grows older but always stays in the background, motivating access to cultural information, reinforcing memory and supporting the creation of meaning [17]. Language development strongly depends on inter-subjective experiences: from the effective engagement of minds and bodies depends cultural learning [9]. Although humans appear to be born with a natural disposition towards cultural learning [21], successful acquisition of cultural skills depend on the interaction quality, especially considering social feedback. Given the social nature of cultural transfer, it is not enough to expose children to new words without providing an adequate *context* to them. Engaging and meaningful activities are especially important to attract interest in the children and *show* them how words can provide the natural pleasure that comes with gaining competence in interacting with their loved ones and with peers. Storytelling has been demonstrated to be a powerful mean to accomplish this as children are born with "[...] an abundant and early armament of narrative tools" [3, p. 90]. Through storytelling, children acquire skills related to the so called *emergent literacy* [18], which is a necessary prerequisite to mastering reading and writing. These skills cover metalinguistic awareness, cohesion and reference in oral communication and the capability of making one's own intentions known to others. Emergent literacy capabilities "[...] are acquired first in language play and in storytelling. Many of them are acquired in the context of childrens interactions with peers, in early play contexts." [5, p. 76]. Once again, the importance of social context and playful interaction is highlighted concerning the acquisition of literacy skills. Wordplay for children appears to be based on matching and substituting words on the basis of their sound rather than their meaning as they appear to [5, p. 78] "[...] derive tremendous pleasure from rhyming words ("you silly"; "no, you pilly") or words that sound similar (adult: "Indians lived in a teepee"; child: "pee-pee!"). In order to become meaningful and precious for children, teaching activities need to have a basis of experiences showing language as a tool to provide pleasure in social activities. In this paper, we will present a software architecture designed to present discrimination tests in a playful setup depicting a social situation with different kinds of virtual agents. This ongoing work builds upon the experience of the Colorado Literacy Tutor [6] and of the Italian Literacy Tutor [7].

SYSTEM ARCHITECTURE

The scripted approach has the disadvantage of not being able to adjust the test depending on the subject's performance. As a limited amount of time is available to administer the test before the child gets tired, choosing the most informative stimulus at each step of the test would represent an advantage when information is clearer on some traits and more uncertain on others. Being able to concentrate on collecting information on specific aspects of linguistic competence that have been observed to be challenging for the child would optimise the

available time. The system architecture we designed to administer the discrimination tests has two main purposes:

- dynamically adapt the test to the child's performance;
- support groups of virtual agents to establish social setups

To pursue the first goal, we represent the discrimination test as a dialogue model where each stimulus, once paired with the child's answer, generates a new stimulus as a system response. This stimulus is selected depending on a utility function taking into account linguistic knowledge and the child's performance. From an architectural point of view, this reflects in a dialogue manager acting as the system's controller and in linguistic knowledge being distributed between the dialogue manager and a database of Italian words. The dialogue manager is provided with the capability to establish which kind of information can be obtained by presenting each available stimulus and with a non-words generator using phonotactic rules to avoid structures not belonging to the Italian language. The database contains morpho-syntactic, phonological and frequency data about words to improve the quality of the selected stimuli. To present the discrimination test in a social setup, the dialogue manager controls a set of virtual agents with different characteristics. In our case, a virtual avatar is presented on a computer screen and acts as the game's guide while a social robot is used to implement a learning-by-teaching approach, detailed in Section . The virtual avatar is controlled using the Unreal Engine 4¹ and its voice is dynamically generated using the Mivoq Voice Synthesis Engine². The synthetic voice has a number of advantages: it allows the system to be easily updated as the proposed stimuli are not pre-recorded, it allows the 3D characters to address the child by calling her by name, thus establishing a closer relationship, and it can be adapted to different kinds of characters. In the specific case of Mivoq, personalised voices and specific prosodic styles can also be synthesised, opening to a number of applications for game-like software artefacts. A tablet interface, also controlled using the Unreal Engine 4, is provided to the child to evaluate the proposed stimuli. Since the ability to adequately use a tablet interface appears to be reliable for 5 years old and onwards children [24], this is the minimum age recommended to apply this technology. The robot used in our implementation is Nao³, which is a well established robotic platform to work with children. The dialogue manager does not make assumptions about the nature of the virtual agents it is connected to. The commands it generates are the same for both the robotic platform and for the 3D character (i.e. *Synthesise, Speak...*). Command implementation is delegated to the specific platform to separate the test logic from its actual implementation. The full schema of the architecture we present is shown in Figure 1. In the following sections, we detail how each module was designed and its role in the general setting. While the system we are developing is able to administer the test without human supervision, we do not exclude the human expert from the experimental setup. The presence of a reference human figure is important to reassure

¹www.unrealengine.com

²www.mivoq.it

³www.softbank.jp/en/corp/group/sbr/

the child and to integrate the obtained results in the light of direct observation of the child's behaviour. In these development stages, moreover, the experience of practitioners is precious to improve the quality of the overall experience without altering the validity of the test.

LINGUISTIC KNOWLEDGE BASE

With the advent of the Big Data and, in particular, with the increasing availability of Linked Open Data, the need to establish a representation format suitable for dynamic, rapidly changing and interconnected *objects* arose. RDF represents the most widely used solution to this need and has been adopted to implement the most widely known repositories of linked knowledge available today. An alternative to RDF is now represented by graph databases. Neo4J [25] is the graph database solution we used in our architecture. It is an open source graph database manager that has been developed over the last sixteen years and has been applied to a high number of tasks related, among others, to data representation [8] and visualisation [11]. In Neo4J, nodes and relationships may be assigned *labels*, which describe the type of the object they are associated to. In this work, labels are mainly used to represent morpho-syntactic characteristics of words and the nature of the relationships among nodes. Nodes and relationships may have *properties*, which are used here to store the details of each single node or relationship. Labels and properties are the main way used by Neo4J to filter data and retrieve answers to user queries. In this work, we use the MultiWordNet-Extended (MWN-E) dataset [14], as the knowledge base to control the decision process for the discrimination test. The MWN-E dataset is based on the MultiWordNet dataset [15] and extended by introducing morpho-syntactic data (e.g. gender, number...), derived forms (e.g. plurals, conjugations...) and SAMPA pronunciations. Also, phonological neighbourhoods are computed and are of particular interest for this work. A word *A* is defined to be a phonological neighbour of the word *B* if it is possible to obtain *B* by altering the phonological representation of *A* using exactly one Insertion/Deletion/Substitution operation. Phonological neighbourhoods are represented by establishing relationships of type HAS_PHONOLOGICAL_NEIGHBOUR between two words if the Minimum Edit Distance of their phonological transcriptions equals 1. This kind of relationship has a *distance* property that, in these cases, is set to 1. Relationships of type HAS_PHONOLOGICAL_NEIGHBOUR are also established between words that have the same pronunciation but have different written forms. In this case the value of the *distance* property is set to 0. Other than the data included in the version presented in [14], the MWN-E version used in this work also contains frequency data for the terms in the vocabulary presented in the *Primo Vocabolario del Bambino* (Children's first vocabulary) [4] and from the Italian Wikipedia⁴. Currently, MWN-E consists of 292282 nodes containing 1536550 properties. 943174 relationships among these nodes are found, phonological neighbourhood relationships at distance 1 representing the majority. The querying language used to extract data from a Neo4J database is Cypher. Cypher

is designed to be a *declarative* language that highlights patterns structure by using an SQL inspired *ascii-art syntax*. A brief overview of the syntactic elements of Cypher queries is given here to help understanding the example queries presented in this paper. The reader is referred to the online Cypher manual⁵ for a more detailed presentation of Cypher. As in graphical representations of graphs nodes are usually represented by circles, in Cypher nodes are represented by round brackets. For example, the query `MATCH (n:VERB) RETURN n` returns all the nodes of the graph labelled as verbs. In the same way, since relationships are usually represented by labelled arrows in graph schemas, relationships between nodes are described by using ASCII *arrows*, too. The query `MATCH (m)-[:DERIVES_FROM]->(:VERB word:'essere') RETURN m` returns all the nodes that contain a term that derives from the *essere* (to be) verb. The SQL-like WHERE clause may also be used to filter results using boolean logic. The query shown in Figure 2 shows how to obtain a pair (w_1, w_2) consisting of dysyllabic words that are phonological neighbours and are obtained by substituting the /p/ phoneme in the first word with the /b/ phoneme in the second word. Sets of words to be excluded after having been presented are also included (in this example, *cubo* and *cupo*) as well as the sorting logic. The first part of the Cypher query matches words that are linked by phonological neighbourhood relationships at distance 1, regardless of arc orientation. A filter is then applied on the syllabic structure using a regular expression on the SAMPA transcription property. In this case, only words presenting a CV-CV structure with the accent on the first syllable and presenting the phonemes /p/ and /b/ in opposition on the head of the second syllable are accepted. The regular expression is dynamically generated by the dialogue system depending on the opposition to present and on the word structure complexity. The former comes from a decision process implemented in the dialogue manager while the latter becomes more complex as the words available for the each considered structure become less informative, as in the case of words presenting oppositions that have already been investigated.

OPENDIAL

Opendial [13] is a dialogue management framework based on probabilistic rules aiming at merging the best of rule-based and probabilistic dialogue management. In cases where a good amount of previous knowledge about the domain is possessed by the dialogue designer with specific needs of fine-tuning rules, the rule based approach can be integrated with probability and utility-based reasoning to fine tune the system's response. Probabilistic rules, in Opendial, are used to setup and update a Bayesian network consisting of variables representing the dialogue state. Depending on this, the dialogue manager selects the most probable user action given a set of, possibly inaccurate, inputs. Using a set of utility functions provided by the dialogue designer, the manager computes the most *useful* system reaction, possibly generating natural language responses or executing actions. In Opendial, it is possible to apply *a priori* estimates on future values of state variables. The probability distributions providing a

⁴Data extracted from the 20/04/2017 Wikipedia.it dump

⁵<https://neo4j.com/developer/cypher-query-language>

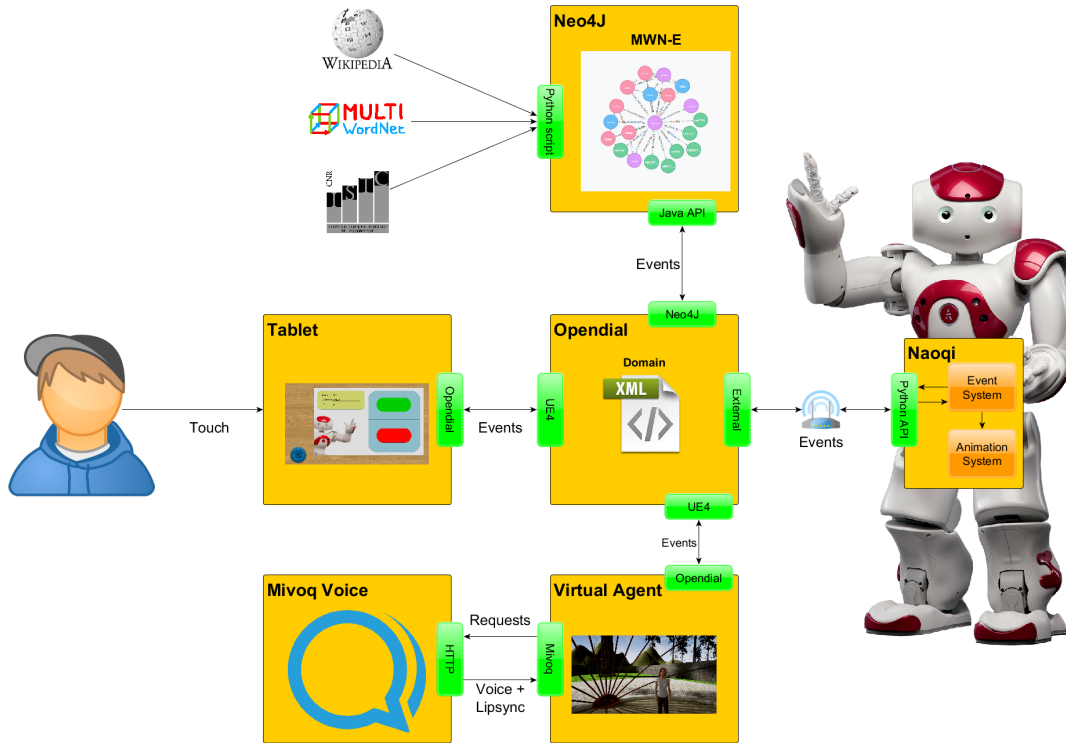


Figure 1. System Architecture.

```

MATCH (n)-[:HAS_PHONOLOGICAL_NEIGHBOUR {distance: '1'}]-(m)
WHERE n.phones =~ '^([p]b|...|LL) (a|e|i|o|u|E|O)1 - (p) (a|e|i|o|u|E|O)$'
AND m.phones =~ '^([p]b|...|LL) (a|e|i|o|u|E|O)1 - (b) (a|e|i|o|u|E|O)$'
AND NOT(n.word IN(['cubo', 'cupo']))
AND NOT(m.word IN(['cubo', 'cupo']))
RETURN DISTINCT n.word AS n, m.word AS m,
CASE
  WHEN n.LEt > 0 AND m.LEt > 0
  THEN 2
  WHEN n.LEt > 0 OR m.LEt > 0
  THEN 1
  ELSE 0
END AS Weight,
(n.LEt + m.LEt)/2 AS PVBMean,
(n.WkFreq + m.WkFreq)/2 AS WkMean
ORDER BY Weight DESC, PVBMean DESC, WkMean DESC
LIMIT 1

```

Figure 2. Example query. Extracts a word pair of disyllabic phonological neighbours opposing the /p/ sound and the /b/ sound in the head of the second syllable.

priori estimates can be updated, using Bayesian inference, after the actual observation arrives to dynamically improve the model. In Opendial, dialogue domains are described in an XML format specifically designed for the dialogue system. This is composed of a set of models triggered by variable updates and containing sets of rules to change the dialogue state. Opendial supports unification in its dialogue specification language so that variables can be included to obtain generic rules. In the example shown in Figure 3, a part of the model that identifies opposing traits given two phonemes is presented. The condition for the considered rule to fire is that the two phonemes in the *opposition* variable are not the same one. If the condition is verified, a custom *HasTraits* function is used to determine if the two phonemes have the *sonorant*

trait. Then, the probability of the set of opposing traits to contain the *sonorant* trait is equal to the XOR of the result obtained by applying the *HasTraits* function on the considered phonemes. Opendial can also be extended with Java-based plugins and functions. In our case, we developed a set of plugins to connect the dialogue system to the Neo4J database and to the remote actors providing the user interface. We also developed the custom function to compute the set of opposed traits given two phonemes and a utility model to select the most informative stimulus at each step. The system makes use of the prediction and feedback mechanism provided by Opendial to build the probability distributions describing the likelihood of a subject to discriminate a specific trait. This is used to select the next stimulus that improve the user model the most, given previous answers. This approach results in an adaptive test. The description of the utility model is beyond the scope of this work so we provide only a brief description of the aspects it takes into account. The model considers the information entropy for each trait, the syllabic structures already used to present the available oppositions, the number of traits opposed in each possible phoneme pair and the intrinsic phoneme complexity evaluated on an acquisitional basis [26]. For all these aspects, a specific utility value is computed. The obtained measures are combined into a utility value that is used to select the best stimulus at each step.

INTERFACE

The interface proposed to the child to mask the discrimination test supports a narrative in which the Nao robot wants to learn how to speak and the 3D character needs the child's help to

```

<rule>
  <case>
    <condition>
      <if var="opposition" relation="contains" value="{X};{Y}"/>
      <not>
        <if var="opposition" relation="contains" value="{X};{X}"/>
      </not>
    </condition>
    <effect prob=" (HasTraits({X},Sonorant) -
      HasTraits({Y},Sonorant)) *
      (HasTraits({X},Sonorant) -
      HasTraits({Y},Sonorant))">
      <set var="OpposedTraits" value="{OpposedTraits}+Sonorant"/>
    </effect>
  </case>
</rule>

```

Figure 3. Example rule to check whether two phonemes have the Sonorant trait in opposition. The HasTraits function has been implemented in Java and exposed to the domain specification language.

teach it. A three-polar setup, shown in Figure 4, is established to involve the child in a socially engaging situation. Through this learning-by-teaching approach, the child is given an authoritative role to avoid making him feel threatened or evaluated. When the system starts, an introductory scenario is presented and the 3D character, shown in Figure 5, introduces itself. The scenario ends with the 3D character asking the child to caress Nao in order to wake it up. This has both the goal of providing the invitation to play and to establish physical contact between Nao and the child. Whether the physical attributes of robots constitute an advantage for acceptability *per se* is still a debated issue. In our work, we attempt to fully exploit the physical presence of the robot by presenting tasks that require the child to physically interact with it. By proposing activities that a 3D character simply cannot be involved into, we attempt to capitalise on the robot's potential to provide a more engaging multisensorial experience. Caressing is also a powerful social mean to build attachment. On the other hand, the high level of control over the 3D character movements allows to efficiently represent its higher competence in the considered setup: differently from Nao, this avatar can move the lips and change its facial expressions, providing effective indications on how to continue playing. An advantage of the presented architecture is that different virtual agents can be combined to build the test upon the various advantages they offer. After a tutorial session where Nao performs a small set of funny behaviours, the child is introduced to the actual test. The dialogue manager selects the most appropriate stimulus and coordinates the two agents so that one presents the first (non-)word and the second presents the second. The child is given one possibility to listen to the stimulus again and is required to provide a same/different feedback using an evaluation card that appears on the tablet. The interface to provide feedback is shown in Figure 6.

CONCLUSIONS AND FUTURE WORK

We have presented the work-in-progress on an architectural setup that has been designed to administer gamified discrimination tests. We interpret the test as a dialogue model between the child and a group of virtual characters controlled by a single artificial intelligence. Instead of providing pre-scripted tests, we propose an approach where the test is dynamically generated. The system is able to exploit a significant amount



Figure 4. The experimental setup.

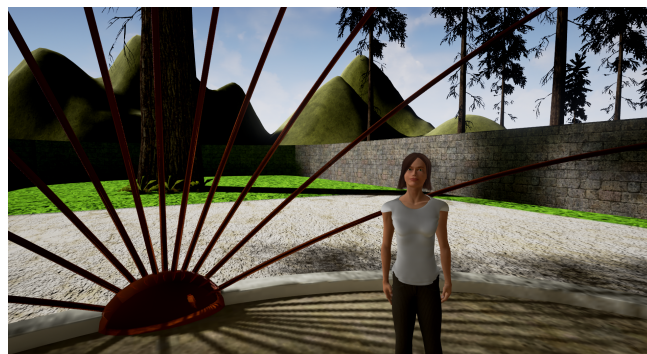


Figure 5. The 3D character. It guides the child through the game and interacts with Nao during cutscenes.

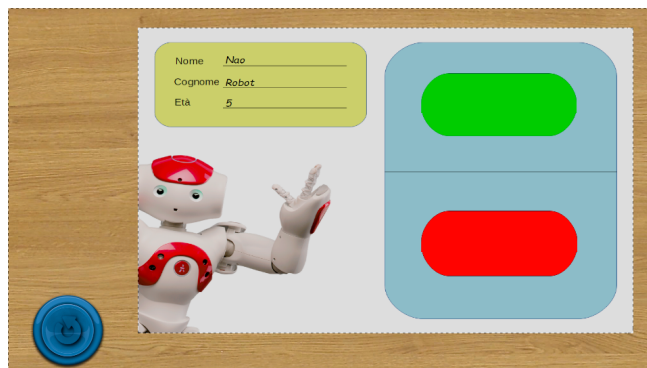


Figure 6. The tablet interface. The child gives feedback by touching the red or green areas to evaluate Nao's performance. A repeat button is also present to allow the child to listen to the opposed words again. This is allowed only once for each stimulus.

of linguistic knowledge to automatically select the most informative stimulus to present at each time. The architecture does not make assumptions about the nature of the virtual agents involved and can be reused to design other types of test. Future work will consist of evaluating the usability and appreciation of the discrimination test we are designing with children that do not show problems in language acquisition to establish a baseline that will be useful to evaluate the approach on children with potential language problems. Also, the possibilities given by the Mivoq engine to train personalised voices will also be explored.

ACKNOWLEDGMENTS

Antonio Origlia's work is supported by Veneto Region and European Social Fund (grant C92C16000250006).

REFERENCES

1. André, C., Ghio, A., Cavé, C., and Teston, B. PERCEVAL: a computer-driven system for experimentation on auditory and visual perception. *CoRR abs/0705.4415* (2007).
2. Brancalioni, A. R., Bertagnolli, A. P. C., Bonini, J. B., Gubiani, M. B., and Keske-Soares, M. The relation between auditory discrimination and phonological disorder. *Jornal da Sociedade Brasileira de Fonoaudiologia* 24, 2 (2012), 157–161.
3. Bruner, J. S. *Acts of meaning*, vol. 3. Harvard University Press, 1990.
4. Caselli, M. C., and Casadio, P. *Il primo vocabolario del bambino*. Milano: Franco Angeli, 1995.
5. Cassell, J. Towards a model of technology and literacy development: Story listening systems. *Journal of Applied Developmental Psychology* 25, 1 (2004), 75–105.
6. Cole, R. A. Roadmaps, journeys and destinations speculations on the future of speech technology research. In *Eighth European Conference on Speech Communication and Technology* (2003).
7. Cosi, P., Delmonte, R., Biscetti, S., Cole, R. A., Pellom, B., and Vuren, S. v. Italian literacy tutor-tools and technologies for individuals with cognitive disabilities. In *INSTIL/ICALL Symposium 2004* (2004).
8. Dietze, F., Karoff, J., Valdez, A. C., Ziefle, M., Greven, C., and Schroeder, U. An open-source object-graph-mapping framework for neo4j and scala: Renesca. In *International Conference on Availability, Reliability, and Security*, Springer (2016), 204–218.
9. Donald, M. *A mind so rare: The evolution of human consciousness*. WW Norton & Company, 2001.
10. Halliday, M. A. K. *Learning How to Mean—Explorations in the Development of Language*. ERIC, 1975.
11. Jiménez, P., Diez, J. V., and Ordieres-Mere, J. Hoshin kanri visualization with neo4j. empowering leaders to operationalize lean structural networks. *Procedia CIRP* 55 (2016), 284–289.
12. Kuhl, P. K. Early language acquisition: cracking the speech code. *Nature reviews neuroscience* 5, 11 (2004), 831–843.
13. Lison, P., and Kennington, C. Opendial: A toolkit for developing spoken dialogue systems with probabilistic rules. *ACL 2016* (2016), 67.
14. Origlia, A., Paci, G., and Cutugno, F. MWN-E: a graph database to merge morpho-syntactic and phonological data for italian. In *Proceedings of Subsidia* (2017).
15. Pianta, E., Bentivogli, L., and Girardi, C. Developing an aligned multilingual database. In *Proc. of the 1st International Conference on Global WordNet* (2002).
16. Polka, L., Jusczyk, P. W., and Rvachew, S. Methods for studying speech perception in infants and children. *Speech perception and linguistic experience: Issues in cross-language research* (1995), 49–89.
17. Reddy, V. *How infants know minds*. Harvard University Press, 2008.
18. Teale, W. H., and Sulzby, E. *Emergent Literacy: Writing and Reading. Writing Research: Multidisciplinary Inquiries into the Nature of Writing Series*. ERIC, 1986.
19. Trevarthen, C. Communication and cooperation in early infancy: A description of primary intersubjectivity. *Before speech: The beginning of interpersonal communication* (1979), 321–347.
20. Trevarthen, C. The functions of emotion in infancy. In *The healing power of emotion: Affective neuroscience, development & clinical practice (Norton Series on Interpersonal Neurobiology)*, D. Fosha and S. M. F. Siegel, D. J., Eds. WW Norton & Company, 2009, 55–85.
21. Trevarthen, C., and Aitken, K. Regulation of brain development and age-related changes in infants. *Motives: The Developmental Function of Regressive Periods*. In M. Heimann (ed.) *Regression Periods in Human Infancy*. Mahwah, NJ: Erlbaum (2003), 107–184.
22. Trevarthen, C., Hubley, P., et al. Secondary intersubjectivity: Confidence, confiding and acts of meaning in the first year. *Action, gesture and symbol: The emergence of language* (1978), 183–229.
23. Tsao, F.-M., Liu, H.-M., and Kuhl, P. K. Speech perception in infancy predicts language development in the second year of life: A longitudinal study. *Child development* 75, 4 (2004), 1067–1084.
24. Vatavu, R.-D., Cramariuc, G., and Schipor, D. M. Touch interaction for children aged 3 to 6 years: Experimental findings and relationship to motor skills. *International Journal of Human-Computer Studies* 74 (2015), 54 – 76.
25. Webber, J. A programmatic introduction to neo4j. In *Proceedings of the 3rd annual conference on Systems, programming, and applications: software for humanity*, ACM (2012), 217–218.
26. Zmarich, C., and Bonifacio, S. Phonetic inventories in italian children aged 18-27 months: a longitudinal study. In *INTERSPEECH* (2005), 757–760.