
Reservoir Computing for Edge-based Automatic Speech Recognition

Nicolo Micheletti - 2024280039
Diego Cerretti - 2024280040
Thomas Adler - 2024389002
29 October 2024

1 Introduction

Automatic Speech Recognition (ASR) ensures seamless interaction between humans and LLM-powered AI. Current state-of-the-art ASR models are transformer-based neural networks [1] that have a very high level of accuracy but come with the cost of high complexity, partly due to the attention mechanisms present in the model [1]. The latest ASR model by Open AI, Whisper Large, has 1.55bn parameters (2.9 GB) [11] and is reported by users ¹ to require around 12 GB of VRAM to run. A model of this size has to be deployed on the cloud, which introduces network latency, slowing response times and degrading user experience. Indeed, edge devices cannot host a model of this size: the latest iPhone 15 Pro Max is estimated to have around 8 GB of RAM. Due to limited resources, current edge-based ASR models also struggle with accuracy [6]. Reservoir Computing offers the potential for a new generation of edge-based ASR models with low latency and high accuracy.

2 Background

Reservoir Computing (RC) [15, 9] is an ideal candidate for an edge-based ASR model because of its low training requirement and strong predictive capabilities in complex time series. The key property of an RC is that its internal layer is fixed, requiring no training. As a result, RCs require very little memory, computational power, fine-tuning or retraining. They can potentially be hosted on memory-constrained edge devices and avoid the network latency issues introduced by cloud-based models. RCs can achieve high-performance levels on specific tasks because of the rich non-linear and recurrent dynamics within the reservoirs. This enables them to extract complex spatio-temporal features when transforming input data into a higher-dimensional space. We further describe the RC architecture in Appendix A.

The following system of equations describes an RC.

$$\begin{cases} x(t+1) = (1 - \gamma)x(t) + \gamma f(Wx(t) + W^{in}u(t) + b), \\ y(t) = W^{out}x(t), \end{cases} \quad (1)$$

It operates with discrete time steps, denoted by t . The non-linear activation function is represented by f , and the reservoir is the internal weight matrix W . W^{in} is the input weight matrix, defining the nodes in the network that receive inputs (input-to-reservoir mapping), while W^{out} represents the output weight matrix (reservoir-to-output mapping). The leakage rate γ controls the amount of past information passed to the next time step. The input at time t is denoted by $u(t)$, and b is a bias term [15, 8].

¹<https://github.com/openai/whisper>

29 3 Related Work

30 **Optimization strategies** State-of-the-art ASR models achieve high accuracy but require complex
31 architectures. Several techniques can be used to alleviate this issue. **Pruning** and **quantization**
32 techniques can be combined to compress a model. Pruning cuts links in a dense network to reduce
33 complexity, whereas quantization lowers the precision of parameters to save memory [7, 4]. Other
34 methods have been implemented to improve convergence rates and training efficiency of ASR
35 models, like **Sortagrad** and **automatic segmentation**. SortaGrad consists of initially training neural
36 networks on shorter audio clips, and gradually increasing their length [2]. Automatic segmentation of
37 input audio into meaningful units enhances the training of automatic speech recognition systems by
38 improving how input features align with phonetic labels and simplifies the data the model needs to
39 handle [5].

40 **ASR for Edge Devices** Recent works have aimed to reduce the latency and memory footprint of
41 ASR models to run on resource-constrained devices. Specifically, Gondi et al. [6] achieved this by
42 building transformer-based models with quantization and other optimization techniques. Xu et al.
43 [14] ran ASR on low-memory devices using Conformer CTC (Connectionist Temporal Classification
44 Automatic Speech Recognition) models. Conformer CTC ASR is a speech recognition model that
45 combines Conformer neural architecture, known for its efficiency in capturing both local and global
46 speech patterns, with CTC loss, which aligns predictions with input sequences without requiring
47 pre-labeled data alignment. Compared to traditional models, these systems exhibit an increase in
48 Word Error Rate (WER), a common accuracy metric in speech recognition tasks.

49 **Reservoir Computing for ASR** Recent works have attempted to reduce complexity while main-
50 taining high accuracy in ASR using RC models. Picco et al. [10] implemented a photonic-based
51 system that makes RC architectures suitable for high-dimensional audio processing tasks. Ansari et
52 al. [3] used heterogeneous single and multi-layer RC models to create non-linear transformations of
53 the inputs, capturing temporal context at different scales.

54 4 Proposal

55 Lighter models operating at the edge face challenges in balancing accuracy, latency, and efficiency. RC
56 provides a promising alternative due to its lower inference latency and reduced memory consumption.
57 In this study, we aim to harness RC to improve the performance of edge-based ASR.

58 For this project, we plan to evaluate our model on the English portion of the following two datasets:

- 59 • The *LibriSpeech*² corpus is a dataset comprising about 1,000 hours of audiobooks from the
60 LibriVox project.
- 61 • The *Common Voice*³ corpus is a multilingual collection of transcribed speech aimed at
62 advancing research in ASR. The dataset contains 9,283 recorded hours.

63 We will compare our implementation against current state-of-the-art approaches in edge-based ASR:
64 Whisper Small from OpenAI [11] and Wav2Vec 2.0 from Meta [6]. Both models performed well
65 on ASR benchmarks but still struggled with accuracy compared to cloud-models, hampering their
66 effectiveness.

67 To evaluate the performance of the models, we will use Word Error Rate (WER), Character Error
68 Rate (CER), and frame-wise accuracy, representing the correct prediction of a sound present in a
69 short audio segment. In addition, the models will be compared based on their memory usage and
70 inference time.

71 To improve the accuracy of our RC model, we will explore various approaches, including those
72 mentioned in Section 3. In addition, we will fine-tune RC-specific parameters. These include reservoir
73 size and connectivity, the number of reservoir layers, recurrent weights, spectral radius, and the
74 degree of chaotic dynamics within the reservoir [12, 16].

²<https://paperswithcode.com/dataset/librispeech>

³<https://paperswithcode.com/dataset/common-voice>

75 References

- 76 [1] Sadeen Alharbi et al. “Automatic Speech Recognition: Systematic Literature Review”. In:
77 *IEEE Access* 9 (2021), pp. 131858–131876. DOI: 10.1109/ACCESS.2021.3112535.
- 78 [2] Dario Amodei et al. “Deep speech 2: End-to-end speech recognition in english and mandarin”.
79 In: *International conference on machine learning*. PMLR. 2016, pp. 173–182.
- 80 [3] Zohreh Ansari, Farzin Pourhoseini, and Fatemeh Hadaeghi. “Heterogeneous reservoir comput-
81 ing models for persian speech recognition”. In: *2022 International Joint Conference on Neural*
82 *Networks (IJCNN)*. IEEE. 2022, pp. 1–7.
- 83 [4] Alexei Baevski et al. “wav2vec 2.0: A framework for self-supervised learning of speech
84 representations”. In: *Advances in neural information processing systems* 33 (2020), pp. 12449–
85 12460.
- 86 [5] Ronan Collobert, Christian Puhersch, and Gabriel Synnaeve. “Wav2letter: an end-to-end
87 convnet-based speech recognition system”. In: *arXiv preprint arXiv:1609.03193* (2016).
- 88 [6] Santosh Gondi and Vineel Pratap. “Performance Evaluation of Offline Speech Recognition
89 on Edge Devices”. In: *Electronics* 10.21 (2021). ISSN: 2079-9292. DOI: 10.3390/
90 electronics10212697. URL: <https://www.mdpi.com/2079-9292/10/21/2697>.
- 91 [7] Song Han et al. “Ese: Efficient speech recognition engine with sparse lstm on fpga”. In:
92 *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate*
93 *Arrays*. 2017, pp. 75–84.
- 94 [8] Mantas Lukoševičius. “A Practical Guide to Applying Echo State Networks”. In: *Neural*
95 *Networks: Tricks of the Trade: Second Edition*. Ed. by Grégoire Montavon, Geneviève B. Orr,
96 and Klaus-Robert Müller. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 659–
97 686. ISBN: 978-3-642-35289-8. DOI: 10.1007/978-3-642-35289-8_36. URL: https://doi.org/10.1007/978-3-642-35289-8_36.
- 99 [9] Mantas Lukoševičius and Herbert Jaeger. “Reservoir computing approaches to recurrent
100 neural network training”. In: *Computer Science Review* 3.3 (2009), pp. 127–149. ISSN: 1574-
101 0137. DOI: <https://doi.org/10.1016/j.cosrev.2009.03.005>. URL: <https://www.sciencedirect.com/science/article/pii/S1574013709000173>.
- 103 [10] Enrico Picco, Alessandro Lupo, and Serge Massar. “Deep photonic reservoir computer for
104 speech recognition”. In: *IEEE Transactions on Neural Networks and Learning Systems* (2024).
- 105 [11] Alec Radford et al. “Robust Speech Recognition via Large-Scale Weak Supervision”. In:
106 *Proceedings of the 40th International Conference on Machine Learning*. Ed. by Andreas Krause
107 et al. Vol. 202. Proceedings of Machine Learning Research. PMLR, July 2023, pp. 28492–
108 28518.
- 109 [12] Bin Ren and Huan-Fei Ma. “Global optimization of hyper-parameters in reservoir comput-
110 ing”. In: *Electronic Research Archive* 30 (May 2022), pp. 2719–2729. DOI: 10.3934/era.
111 2022139.
- 112 [13] David Verstraeten et al. “The unified reservoir computing concept and its digital hardware
113 implementations”. In: *Proceedings of the 2006 EPFL LATSIS Symposium*. 2006, pp. 139–140.
- 114 [14] Mingbin Xu et al. “Conformer-based speech recognition on extreme edge-computing devices”.
115 In: *arXiv preprint arXiv:2312.10359* (2023).
- 116 [15] Min Yan et al. “Emerging opportunities and challenges for the future of reservoir computing”.
117 In: *Nature Communications* 15.1 (Mar. 2024), p. 2056.
- 118 [16] Bolin Zhao. “Seeking optimal parameters for achieving a lightweight reservoir computing: A
119 computational endeavor”. In: *Electronic Research Archive* 30 (June 2022), pp. 3004–3018.
120 DOI: 10.3934/era.2022152.

121 **A Reservoir Computing Architecture**

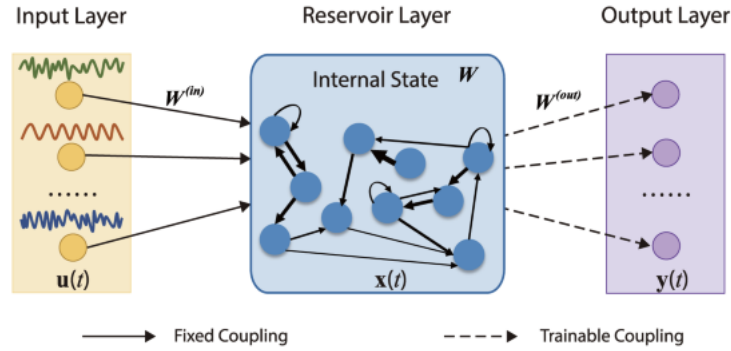


Figure 1: RC architecture [15] (Authorized under (CC BY 4.0))

122 As Figure 1 shows, an RC consists of three layers: an input (sensing) layer, a reservoir (processing)
 123 layer, and an output (control) layer. The input layer sends data to the reservoir, a fixed-weighted
 124 network that projects input data into a higher-dimensional feature space. The output layer, the only
 125 trainable component, typically uses linear regression to map these signals to the final output. The
 126 reservoir is initialized once, with its size, connectivity, and chaotic dynamics fixed. Then, the RC
 127 processes sequential data at every time step, and its recurrent connections ensure past data is carried
 128 over to future time steps [15, 13, 9, 8].