

Radar based Gesture Recognition

Lakshya Kumar Tangri 2018EE10222

Somanshu Singla 2018EE10314



Faculty Supervisor: Prof. Seshan Srirangarajan

17 December 2021

Introduction



- Human gestures: a more natural interface between humans and computers.
- Cameras: Alone or together with other devices raise major privacy concerns

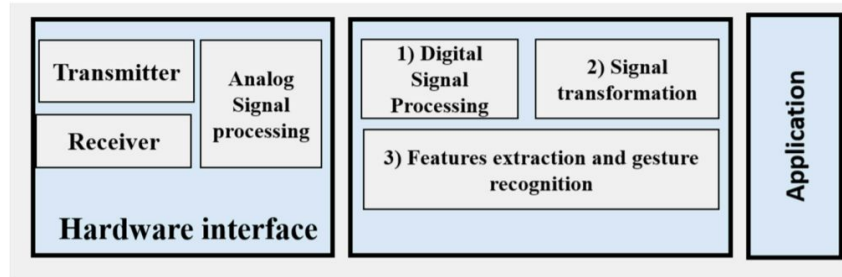
Short-range radars have the ability to detect micro-movements with high precision and accuracy making it a good candidate for realizing gesture-based interfaces

Source: Jaime Lien, Nicholas Gillian, M. Emre Karagozler, Patrick Amihoud, Carsten Schwesig, Erik Olson, Hakim Raja, and Ivan Poupyrev. 2016. Soli: ubiquitous gesture sensing with millimeter wave radar. ACM Trans. Graph. 35, 4, Article 142 (July 2016), 19 pages.
DOI:<https://doi.org/10.1145/2897824.2925953>

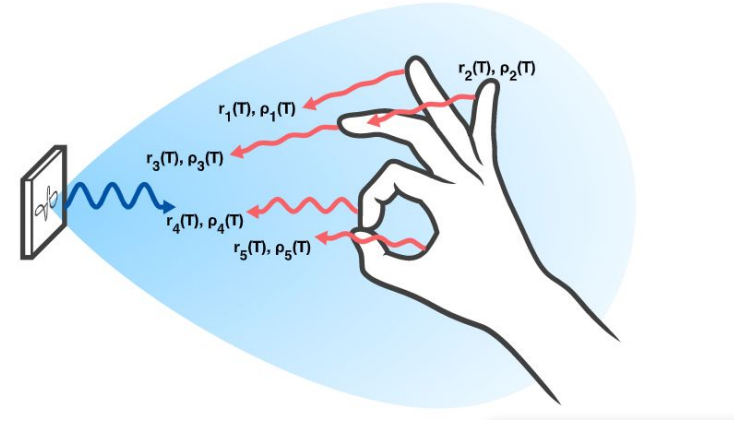
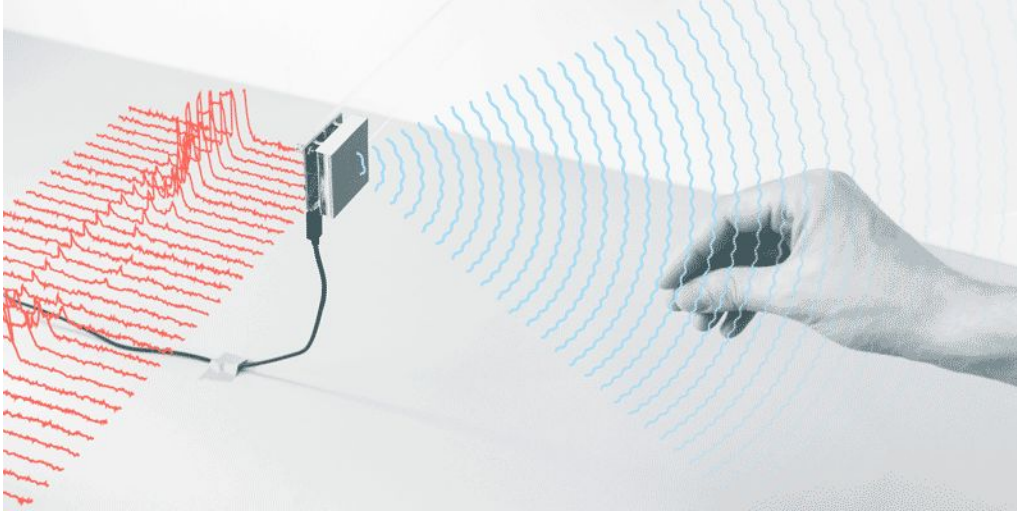
Project Objective

To create a robust, radar based power efficient system which can classify predefined gestures accurately in real time on a mobile device

- **Robustness:** With respect to different users(inter-user) and different gesture instances of a single user(inter-session)
- **Accuracy:** To be able to correctly classify the input gesture
- **Power Efficiency and Real Time:** Ability to recognize gesture quickly for a seamless user experience without draining too much battery

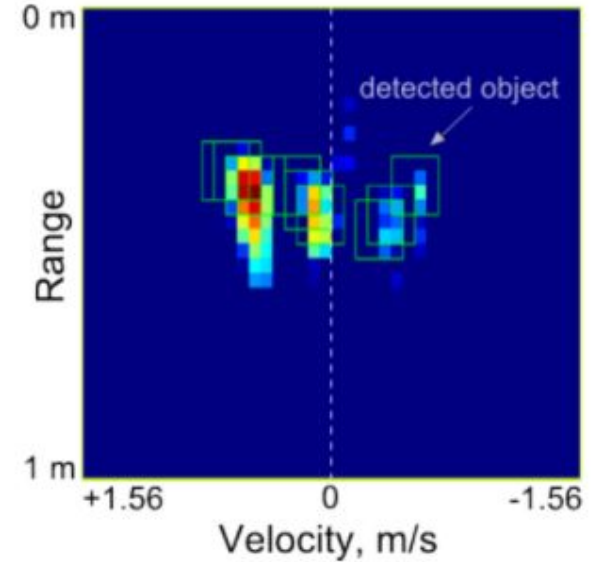
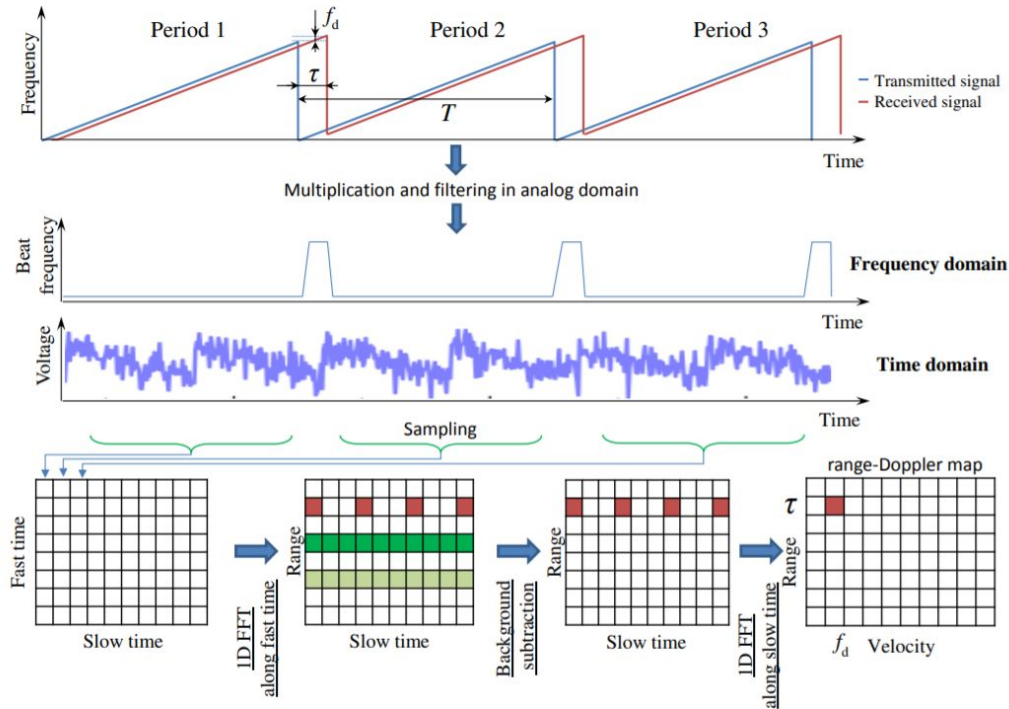


Working



Source: Jaime Lien, Nicholas Gillian, M. Emre Karagozler, Patrick Amihoud, Carsten Schwesig, Erik Olson, Hakim Raja, and Ivan Poupyrev. 2016. Soli: ubiquitous gesture sensing with millimeter wave radar. ACM Trans. Graph. 35, 4, Article 142 (July 2016), 19 pages.
DOI:<https://doi.org/10.1145/2897824.2925953>

Range Doppler



Source: P. Molchanov, S. Gupta, K. Kim and K. Pulli, "Short-range FMCW monopulse radar for hand-gesture sensing," 2015 IEEE Radar Conference (RadarCon), 2015, pp. 1491-1496, doi: 10.1109/RADAR.2015.7131232.

Literature Survey

Types of papers:

- I. Analyzing properties of FMCW Radar
- II. Traditional machine learning
- III. Deep Learning

Observations:

- Majority of the literature leave out details such as memory, power consumption, model robustness, prediction time
- No standardised dataset large/diverse enough to be used as a common benchmark for testing these models
- The variable sequence length of input data has not been emphasized

Experiments

Experimented with CNN architecture using variety of features.
Used sub sampling and padding techniques to handle variation in length of inputs received.

Model Architecture: Modified LeNet⁷

Layer (type)	Output Shape	Param #
Conv2d-1	[-1, 6, 32, 32]	6,006
MaxPool2d-2	[-1, 6, 16, 16]	0
Conv2d-3	[-1, 12, 12, 12]	1,812
MaxPool2d-4	[-1, 12, 6, 6]	0
Linear-5	[-1, 216]	93,528
Linear-6	[-1, 54]	11,718
Linear-7	[-1, 11]	605
=====		
Total params: 113,669		
Trainable params: 113,669		
Non-trainable params: 0		

Background

$$RP_t^i = \sum_v RD^i(r, v, t)$$

$$DP_t^i = \sum_r RD^i(r, v, t)$$

Here,

RD is the range Doppler

RP is the range profile

DP is the doppler profile

i is receiver number or channel

t is the frame number

r is the range bin

v is the velocity bin

Motion Profile is a concatenation of Range Profile and Doppler Profile

Experiment Results

Compared accuracy metrics of some basic features:

All the accuracy metrics reported are for 50-50 Train Test Split.

By using Range profile as feature:

Accuracy = 84.897%

By using Doppler profile as feature:

Accuracy = 89.589%

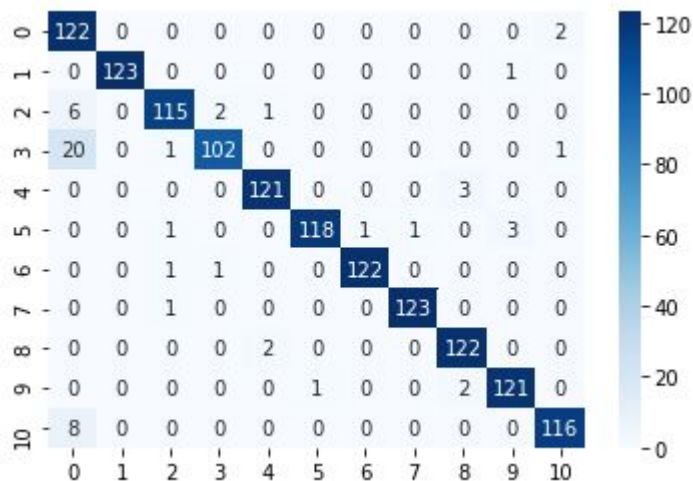
By using Motion Profile feature:

Accuracy = 95.45%

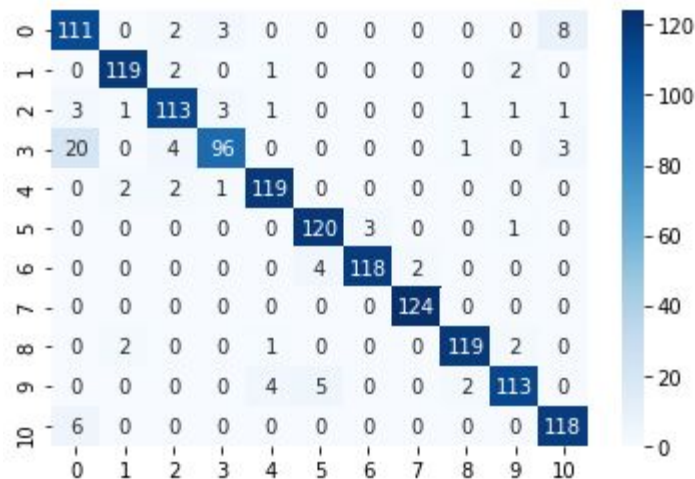
By using complete sequence of RDIs:

Accuracy = 93.40%

Confusion Matrix on Unpruned Data Set



For Motion Profile feature



For Sequence of RDIs

Gesture Set Pruning

- Observed that features mentioned in the literature don't do well on the complete gesture set and after viewing the confusion matrices concluded that out of all gestures, model's performance is not upto the mark for some gestures(i.e it confuses between two similar gestures) and good for others.
- In a practical setting we don't really need a gesture set containing 10- 15 gestures even some 5 - 7 gestures (being ergonomic) can suffice , so keeping this in mind we test our model again on removing of these bad performing gestures



Accuracy Metrics

Conducted robustness tests such as inter user testing, inter session testing
Pruned the data set based on the features which were not confusing and
repeated the same tests (we removed gestures 3 and 10).

The results are as follows:

Feature	Baseline	Leave One User Out (LOUO)	Leave One Session Out (LOSO)	Baseline [Pruned Data]	LOUO [Pruned Data]	LOSO [Pruned Data]
Motion Profile	95.67%	90.18%	92.84%	98.12%	93.87%	97.95%
Sequence of RDIs	93.11%	86.22%	90.33%	96.86%	91.95%	96.58%

CNN Model Complexity

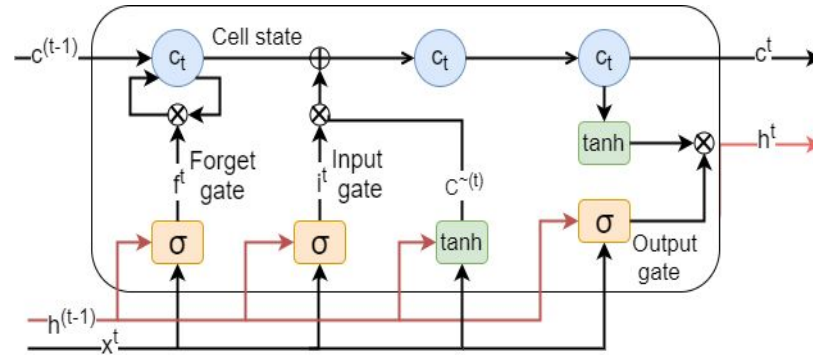
These statistics can help us know about the model complexity:
FLOP denotes Floating Point Operations during the inference stage.

Model	Time taken (ms)	FLOP (M)	Number of parameters
Motion Profile	0.89	3.91	766547
RDI	1.90	13.02	113669

- These are calculated for a 40-frame gesture during the inference stage.
- The time taken is calculated on an i5, 1.8GHz machine.

LSTM Encoder

- LSTM encoder is a model which can take a variable length sequence data and can return a fixed shape representation of our input.
- This representation of the input is encapsulated in “hidden state” of the model.
- Hidden state is basically a vector which is comprised of values from each of the LSTM cells, we get a hidden state from the model at every timestep.
- Typically the hidden state at the final time step of LSTM is used as the representation of input data



Attention Mechanism

- Usage of the last hidden state as representation is not optimal as the model tends to forget some information over time
- Attention is a technique to get the representation of input by computing a weighted sum of states at all times where those weights are also learnable

Formally:

$$q^j = \max_{t \in (1 \dots n)} (h_t^j) \quad \hat{h}_t = \frac{h_t}{||h_t||}$$

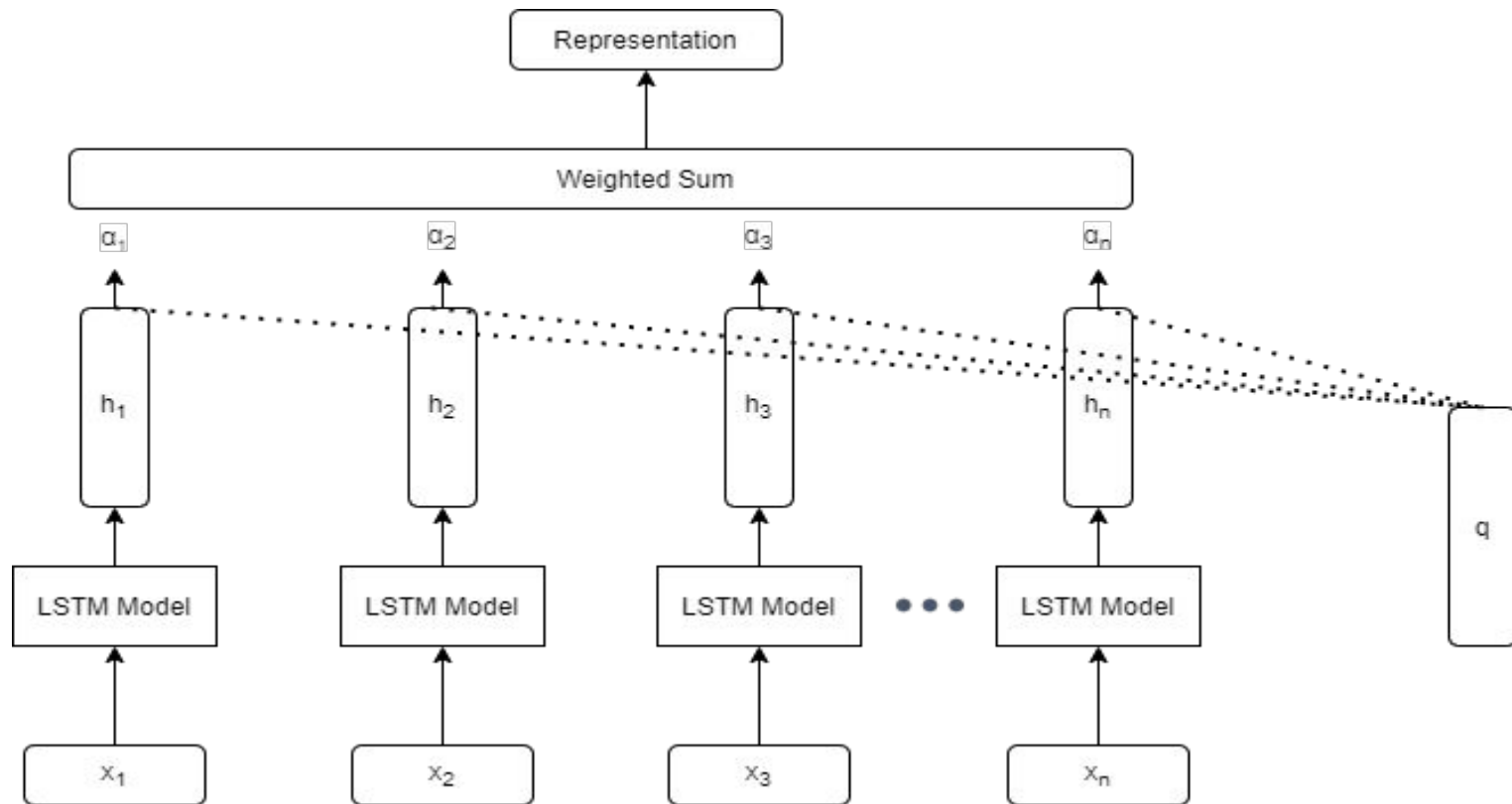
Here,

$$\alpha_t = \frac{\exp(\hat{h}_t^T q)}{\sum_{j=1}^n \exp(\hat{h}_j^T q)}$$

j denotes the jth dimension of h_t
t denotes the tth time step
 h_t denotes the hidden state at tth time step

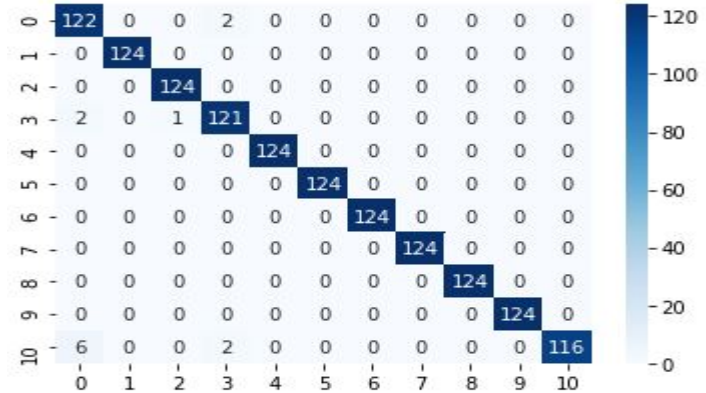
$$Representation = \sum_{t=1}^n \alpha_t h_t$$

LSTM Architecture



LSTM Results

As we can see for this model the two bad performing gestures are also getting classified correctly around 94% of times which decreases the importance of pruning.



Confusion Matrix (2-layer model)

Model	Baseline	Leave One User Out (LOUO)	Leave One Session Out (LOSO)
2-layered	99.04%	92.07%	96.25%
1-layered	98.97%	92.90%	95.92%

LSTM Model Complexity

These statistics can help us know about the model complexity:
FLOP denotes Floating Point Operations during the inference stage.

Model	Time taken (ms)	FLOP (M)	Number of parameters
1-layered LSTM	4.91	15.82	202443
2-layered LSTM	9.29	31.58	334539

- These are calculated for a 40-frame gesture during the inference stage.
- The time taken is calculated on an i5, 1.8GHz machine.

Model	Baseline	LOUO	LOSO	Time taken(ms)	Number of parameters	FLOP (M)
GVLAD ⁷	98.24	91.38	97.75	4.8*	>10M	—
Res3DTENet ⁸	96.99	92.25	—	—	>10M	—
LSTM 2-layered	99.04%	92.07%	96.25%	9.29	334539	31.58
LSTM 1-layered	98.97%	92.90%	95.92%	4.91	202443	15.82
CNN Motion Profile	95.67%	90.18%	92.84%	0.89	766547	3.91
	98.12%	93.87%	97.95%			
CNN Sequence of RDIs	93.11%	86.22%	90.33%	1.90	113669	13.02
	96.86%	91.95%	96.58%			

* time taken is calculated on an i7, 2.2 GHz machine.

Progress Made

- Experimented with various features and architectures proposed in the literature using the publicly available Google Soli Data Set⁶
- Realised that even though the input data can have variable length samples, the current work in this area doesn't really delve deep and instead uses stopgap approaches.
- Experimented with attention based sequence architecture to handle variable length samples.

[6] Wang, S.; Song, J.; Lien, J.; Poupyrev, I.; Hilliges, O. Interacting with soli: Exploring fine-grained dynamic gesture recognition in the radio-frequency spectrum. In Proceedings of the 29th Annual Symposium on User Interface

Future Targets

- Data set collection on TI provided radars and training cum testing of models on dataset for selection of the best performing model
- Implementing the gesture recognition pipeline on actual hardware
- Conducting real-time, power and robustness tests [such as including left handed users, varying hand sizes]

THANK YOU

Questions ?