

**Language models in recognition applications:
a new approach**
**பல்வகை உணர்விகளில் மொழி மாதிரியங்களின்
பயன்பாடு : ஒரு புதிய கண்ணோட்டம்**

A G Ramakrishnan and P Arulmozhi
ஆ. க. ராமகிருஷ்ணன், ப. அருள்மொழி

MILE Laboratory, Department of Electrical Engineering, Indian Institute of Science,
Bangalore
மநுமொழி ஆய்வகம், மின் பொறியியல் துறை, இந்திய அறிவியல் பயிலகம்,
பெங்களூரு.

கட்டுரைச் சுருக்கம் (Abstract)

இந்தக் கட்டுரையானது தமிழில் வாய்மொழி மற்றும் ஒளிவழி எழுத்துரு உணர்தல் (speech and optical character recognition) முறைமைகளில் மொழி மாதிரியங்களின் (language models) பயன்பாட்டில் எதிர்கொள்ளப்படும் சில குறிப்பிட்ட இடர்பாடுகளை விவாதிக்கிறது. இத்தகைய மொழி மாதிரியங்கள் ஐரோப்பிய மொழிகளில் மிகுந்த இலாபகரமாகப் பயன்படுத்தப்பட்டாலும், அவை அவ்வாறே தமிழில் பயன்படுத்தப்படும்பொழுது, தமிழின் மொழியியல் செறிவு காரணமாக அம்மொழிகளில் காணப்படாத சில புதிய இடர்பாடுகளை எதிர்கொள்ள நேரிடுகிறது. இது தமிழின் புணர்ச்சி விதிகள் (sandhi rules) மற்றும் உருபனியல் மாற்றச் செறிவு (morphological richness) காரணமாகவும் சொற்களோடு விசுவாசம் இணைந்து பொருள்தரும் பண்பு (agglutination) மற்றும் பகுதி கட்டில்லாச் சொல்வரிசை முறை (Partially free word order) காரணமாகவும் ஏற்படுகிறது. ஒரு மீப்பெரு பனுவலை (large corpus) அல்லது பனுவலில் இருந்து பெறப்பட்ட பெரும் சொல்வங்கியை n-கிராம் மொழி மாதிரியம் கொண்டு பகுப்பாய்வு செய்ததன் மூலமாகப் பெறப்பட்ட பல்வேறு புள்ளி விவரங்களை முன் வைத்து, அதன் மூலம் தமிழில் எத்தகைய மொழி மாதிரியங்களை எவ்வாறு பயன்படுத்தலாம் என்பது குறித்த எங்கள் கருத்துக்களையும் இங்கு சமர்ப்பிக்கிறோம்.

முன்னுரை (Introduction)

தற்போதைய கணினிமுறை பெருஞ் சொற்களஞ்சிய தொடர் வாய்மொழி உணர்தல் (large vocabulary continuous speech recognition - LVCSR) முறைமைகள் 2,00,000-க்கும் மேற்பட்ட எண்ணிக்கையிலான பெரும் சொல்லகராதியைப் பயன்படுத்துகின்றன. சொற்களஞ்சியத்தின் பருமன் அதிகரிக்கும்போது அதையொத்து விட்டர்பி சொல் அணிக்கோவை (word lattice), இன்ன பிற தேடுதல்களில் கணிப்பு சிக்கற்பாடு அதிகரிக்கின்றது. மொழியியல் செறிவு காரணமாகத் தமிழ் சொல்வங்கியானது பனுவலின் அதிகரிப்பு விகிதத்துக்கு ஏற்ப (மடக்கை விகிதத்தில்) அதிகரிக்கிறது. ஒரு வினை வேர்ச்சொல்லின் (root verb) வருவித்த படிமங்கள் (derived forms) மிகச் சிலவே

உள்ள ஆங்கிலம் மற்றும் ஹிந்தி போன்ற மொழிகளைப் போலல்லாமல் தமிழானது ஒரு வேர்ச்சொல்லுக்குச் சில ஆயிரங்களுக்கு மேற்பட்ட வருவித்த படிமங்களைக் கொண்டுள்ளது [1]. வினைச்சொற்களின் ஒவ்வொரு உருபனியல் கூட்டப்பட்ட வடிவமும் காலம், இடம், பொருள், ஏவல், மறுப்பு, அழுத்தம், கேள்வி போன்ற பல்வேறு வகையான பொருளை வெளிப்படுத்தும் இயல்பில் புதிய எல்லைகளைக் கொண்டுள்ளது. சொல்லிலக்கண விதிகளில் அருகாமை சொல்லிலக்கணக் கூறுகளுடன் இணையும் தன்மையானது ஒருவிதமான நெகிழ்வுத் தன்மையைக் கொண்டுள்ளது. இவ்விதமான நெகிழ்வுத் தன்மையானது சொல்லிலக்கணக் கூறுகள் வெவ்வேறு இடங்களில் இணைய அனுமதிக்கிறது. இதனால் தமிழில் உருபனியல் மாற்றச் செறிவானது வினைச்சொற்களின் பல்வேறு வடிவங்களின் எண்ணிக்கையை ஆச்சரியப்படத்தக்க வகையில் அதிகரிக்கிறது [2]. உதாரணமாக நாங்கள் பயன்படுத்தும் வரையறுக்கப்பட்ட சொல்வங்கியில் உள்ள 'வா' என்கிற வேர்ச்சொல்லின் வருவித்த படிவங்களின் தனிப்பட்ட எண்ணிக்கை 4567 ஆகும். இதன் மூலமாகத் தமிழ் போன்ற செறிந்த மொழிகளுக்கான உணர்தல் மூலோபாயங்கள் (recognition strategies) ஆங்கிலம் மற்றும் ஹிந்தி போன்ற மொழிகளுக்கான மூலோபாயங்களை ஒத்திருக்கவியலாது என்பது புலனாகிறது. எனவே தமிழுக்கு ஏற்றார்போல மாற்று மூலோபாயங்கள் வகுக்கப்பட வேண்டியதின் அவசியம் இங்கு வலியுறுத்தப்படுகிறது.

வாய்மொழி உணர்வியில் மொழி மாதிரியங்கள் (Language models in ASR)

LVCSR முறைமைகள் இரண்டு லட்சம் சொற்களுக்கும் மேற்பட்ட எண்ணிக்கையிலான பெரும் சொல்லகராதியைப் பயன்படுத்துகின்றன. ஆனால் ஒரு சராசரி மனிதனின் மூளையில் பதிந்திருக்கும் சொற்களின் எண்ணிக்கை அதாவது மனிதனின் கற்றுக்கொண்ட சொல்லகராதி இருபதாயிரத்திலிருந்து [3] ஐம்பதாயிரம் [4] வரை தான். இதனால், கணினிகளுக்கு மட்டுமே சாத்தியமாகக்கூடிய ஒரு மீப்பெரு கட்டற்ற சொல்லகராதியானது, உணர்தல் முறைமைகளில் மேலும் குழப்பங்களைத் தருவிக்கவே வாய்ப்பு அதிகம் என்பது அறியப்படுகிறது. நாம் உரையாடும்போது, ஆள்களம் நமக்கு பேசுவதற்கு முன்பே பிடிபட்டுவிடுகிறது. அதன்பின், பல்பொருள் ஒருமொழிச் சொற்கள் பயன்படுத்தப்படும், மனித மூளையின் தேடுதலானது, அச்சமயத்தில் விவாதிக்கப்படும் ஆள்களத்துக்கு உட்பட்டே அமைகிறது. மேலும் கட்டயடுதல் செய்பணி (dictation tasks) முறைகள் பொதுவாக ஒரே அல்லது ஒருசில தொடர்புடைய ஆள்களங்களுக்கு உட்பட்டே அமைகின்றன.

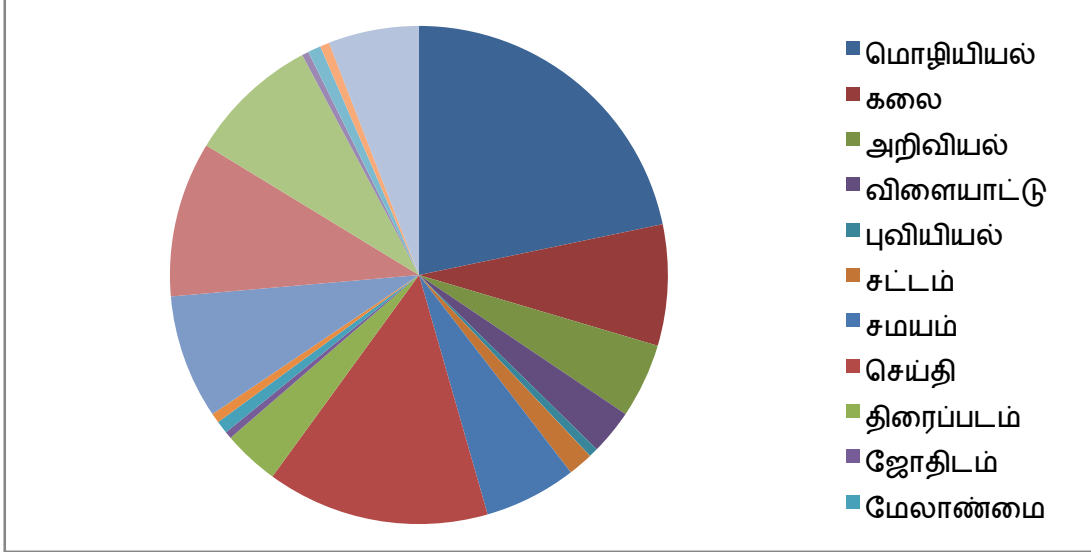
இந்நிலையில் தகுந்த ஆள்களத்தை முன்னறிவதன் மூலமாக அறியப்படும் அறிவுப்புலத்தை வாய்மொழி மற்றும் எழுத்துரு கண்டறிதலில் பயன்படுத்துவது இத்தகைய குழப்பங்களைக் குறைப்பது மட்டுமில்லாது, கணிப்பு சிக்கற்பாட்டையும் குறைத்து துல்லிய அளவையின் படித்தரத்தை உயர்த்தும் என்ற கோட்பாட்டை இங்கு முன்மொழிகிறோம். இந்நோக்கில் இங்கு 2.2 கோடிக்கும் மேற்பட்ட சொற்களைக் கொண்ட பனுவல் பகுப்பாய்வுக்கு உட்படுத்தப்பட்டுள்ளது. இப்பனுவலானது 14 தனிப்பட்ட ஆள்களங்களாக (domains) வகைப்படுத்தப்பட்டுள்ளது. அவை முறையே மொழியியல், கலை, விளையாட்டு, புவியியல், சட்டம், சமயம், அறிவியல், செய்தி, அரசியல், திரைப்படம், ஜோதிடம், மேலாண்மை, பொருளாதாரம் மற்றும் சமூகவியல்

அட்டவணை 1: தமிழ்ச் சொல்வங்கியின் அளவு, தனிச்சொற்களின் எண்ணிக்கை, ஆள்களத்தைத் தேர்வு செய்வதால் பெறப்படும் பேரகராதி மற்றும் சிக்கற்பாட்டு குறுக்கம், வரையறுக்கப்படாத ஒற்றைப் பேரகராதியின் தனிச்சொற்களின் எண்ணிக்கை : ஓர் ஒப்பீடு.

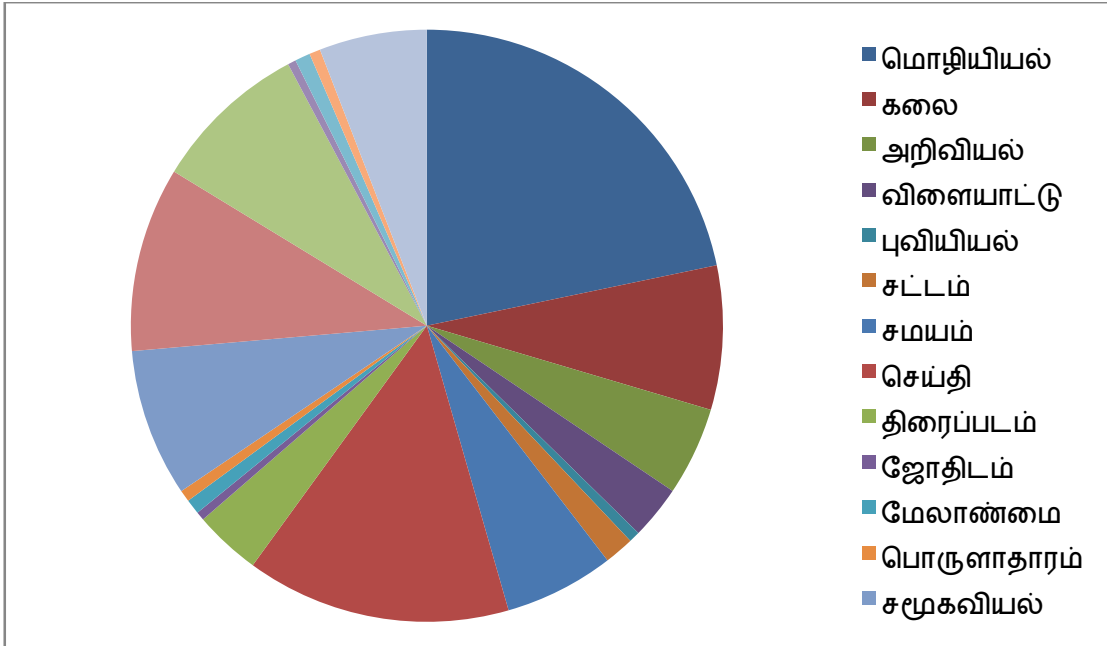
ஆள்களம்	சொற்களின் எண்ணிக்கை	தனிச் சொற்களின் எண்ணிக்கை	சொல்லகராதியின் குறுக்க விகிதம் (%)	ஒப்பீட்டு விகிதம் (ஆள்களத்துக்கு உட்பட்டது - %)
மொழியியல்	4179938	614184	35.8	85.3
கலை	1055330	222333	12.9	78.9
அறிவியல்	769531	136931	7.9	82.2
விளையாட்டு	831656	82384	4.8	90.0
புவியியல்	59124	17051	0.99	71.2
சட்டம்	189321	45738	2.7	75.8
சமயம்	878975	169277	9.8	80.7
செய்தி	6056139	408833	23.8	93.2
திரைப்படம்	779432	103163	6.0	86.7
ஜோதிடம்	37757	13194	0.7	65.1
மேலாண்மை	85183	23116	1.3	72.8
பொருளாதாரம்	120157	17873	1.0	85.1
சமூகவியல்	1309273	227952	13.3	82.6
அரசியல்	3779525	285040	16.6	92.5
சரித்திரம்	1681028	241082	14.0	85.6
வேளாண்மை	42754	12396	0.7	71.0
கல்வி	126546	23749	1.4	81.2
கவிதை	30105	17065	0.99	56.7
மற்ற பிற	699317	166283	9.7	76.2
மொத்தம்	22707568	1716256	--	92.5

ஆகியவை. இவற்றிலிருந்து பல்வகைப்பட்ட ஆள்களங்களில் உள்ள தனிச்சொற்கள், சொல் மற்றும் எழுத்து மட்ட n-கிராம் புள்ளிவிவரங்கள் போன்ற பலவகைப்பட்ட பகுப்பாய்வு முடிவுகள் தருவிக்கப்பட்டுள்ளன. அட்டவணை 1, சில பொதுவான ஆள்களங்களையும் அவை ஒவ்வொன்றிலும் உள்ள தனிச்சொற்களின் எண்ணிக்கையையும் அவற்றின் குறுக்க மற்றும் ஒப்பீட்டு விகிதங்களையும் காட்டுகிறது. உதாரணமாக, மொத்தச் சொல்வங்கியில் உள்ள தனிச்சொற்களின் எண்ணிக்கை 17.16 லட்சமாக இருந்தபோதிலும் 'அறிவியல்' ஆள்களத்தில் உள்ள தனிச்சொற்களின் எண்ணிக்கை 136931 ஆகும். ஆக, நாம் அறிவியல் துறை சார்ந்த ஒரு சொல்வகையின் பணியில் ஈடுபடும்போது, 136931 சொற்களாலான சொல்வங்கியைப் பயன்படுத்தினால் போதுமானது. இது மென்பொருளின் வேகத்தையும் துல்லியத்தையும் நன்கு அதிகரிக்கும் என்பது

திண்ணம். வட்டப் படங்கள் 1 மற்றும் 2-ம் மேற்கூறிய கருத்துக்களை வலியுறுத்துகின்றன.



படம் 1. மநுமொழி தமிழ் பனுவலில் ஆள்களம் வாரியாக மொத்த சொற்களின் எண்ணிக்கை - ஒரு வட்டப் படம்.



படம் 2. மநுமொழி தமிழ் பனுவலில் ஆள்களம் வாரியாக தனிச் சொற்களின் எண்ணிக்கை - ஒரு வட்டப் படம்.

எவ்விரு ஆள்களங்களுக்குள்ளும் பொதுவான சொற்கள் குறைவே (Words common across domains are few)

இரண்டாம் அட்டவணை, ஒரு உதாரணமாக, "சமயம்" என்ற ஒரு ஆள்களத்துடன் வெவ்வேறு ஆள்களங்கள் தனித்தனியே சேரும்போது மொத்தச்சொற்கள் எத்தனை, அவை இரண்டிற்கும் பொதுவான சொற்கள் எத்தனை என்பதைக் காட்டுகிறது. "சமயம்" தனிச்சொற்கள் 169277 உள்ளன. மொத்த சொற்களுடன் ஒப்பிடும்போது பொதுச் சொற்கள் 13 அல்லது அதற்கும் குறைவான விழுக்காடுகளே உள்ளன. ஆகையால் ஆள்களம் பொருத்தே நாம் சொல் வங்கியைப் பயன்படுத்த வேண்டும். அட்டவணை 2 ஒரு மாதிரியே. உண்மையில் 19 ஆள்களங்களுக்கிடையே, 171 இருபுறம் பிணைப்புகளைச் சோதிக்க முடியும்.

அட்டவணை 2: இரு ஆள்களங்களுக்கிடையே உள்ள பொதுவான தனிச்சொற்கள். எல்லா வரிசைகளுக்கும் பொதுவான ஆள்களமான "சமயம்", 169277 தனிச்சொற்கள் கொண்டுள்ளது.

ஆள்களம்	தனிச் சொற்கள்	மொத்தம் (கூட்டு)	பொதுவான சொற்கள்
திரைப்படம்	103163	260853	11587
கலை	222333	345710	45900
விளையாட்டு	82384	240884	10777

பல ஆள்களங்களுக்குப் பொதுவான தனிச்சொற்களின் எண்ணிக்கை (Unique words common across domains)

அட்டவணை 3-ஆனது ஆள்களங்களுக்குப் பொதுவான தனிச்சொற்களின் எண்ணிக்கை, ஆள்களங்கள் சேர்ச் சேர வெகு விரைவில் குறைந்து கொண்டே வருவதைக் காட்டுகிறது. இது, தனிச்சொற்களும் ஒவ்வொரு துறைக்கும் பரத்யேகமானவையே என்பதை உணர்த்துகிறது. ஆக எல்லாத்துறைகளையும் சேர்த்த மொத்தச்சொல் வங்கியை உபயோகிப்பதில் எவ்வித லாபமும் இல்லை என்பது புலனாகிறது.

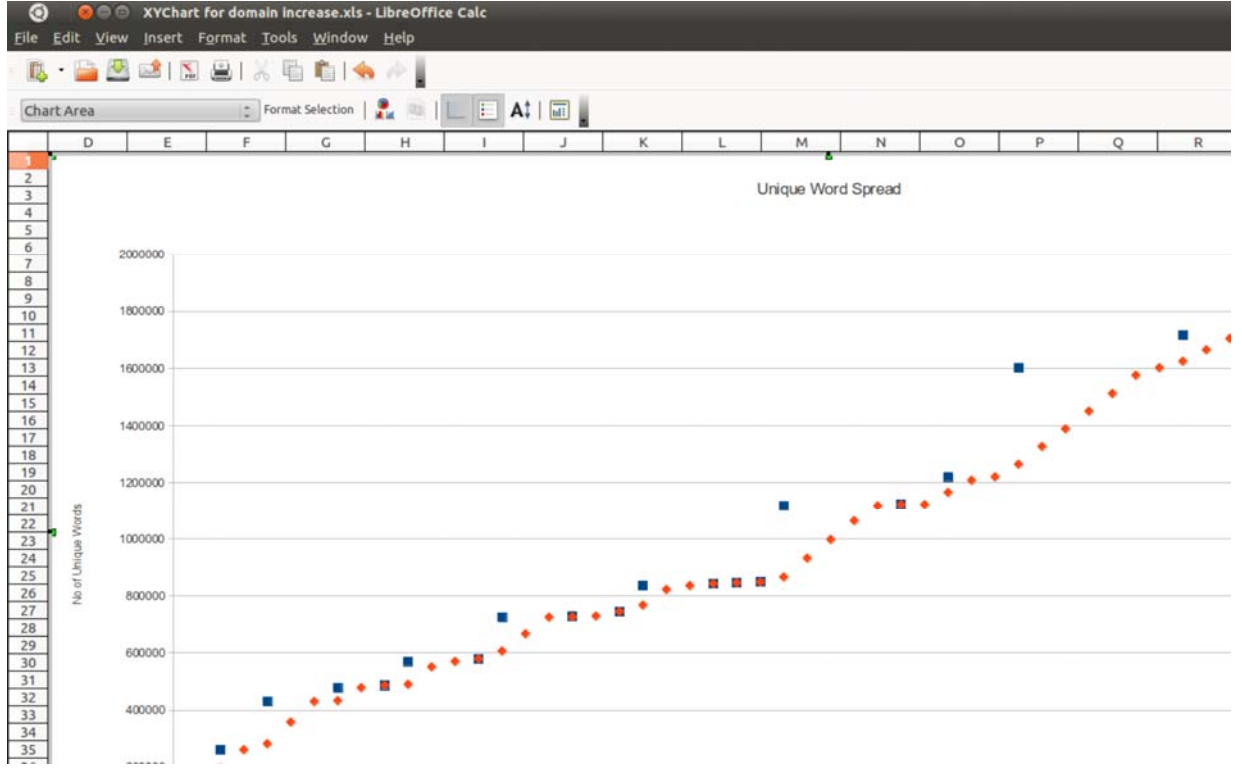
கையெழுத்து உணர்வியில் மொழி மாதிரியங்கள் (Language models in handwriting recognition)

கையெழுத்து மற்றும் அச்சிடப்பட்ட பனுவல்களைக் கையாளும் ஐரோப்பிய மொழிகளுக்கான ஆவணக் கண்டறிகை முறைமைகள் (document recognition systems) கண்டறிதலிலும், உணர்பிழைகளைத் திருத்துவதிலும் மீப்பெரும் சொல்லகராதிகளைப் பயன்படுத்துகின்றன. படம் 3, தமிழ் மொழியின் ஒரு முக்கியமான இயல்பைத் தெளிவாகக் காட்டுகிறது. ஒவ்வொரு ஆள்களமாகச் சேர்த்துக் கொண்டு செல்கையில், தனிச்சொற்களின் மொத்த எண்ணிக்கை, தவிட்டாமல், மென்மேலும் வளர்ந்து கொண்டே இருப்பதுதான் அது. ஆகவே, தமிழ் மொழியானது, ஒரு வரையறுக்க முடியாத தனிச்சொல் வங்கியை உடையது. அதனால், சொல் மட்டப் பேரகராதி (word level lexicon) கொண்டு பிழை திருத்துவது

அட்டவணை 3: இரண்டு, மூன்று, நான்கு எனப் பல ஆள்களங்களுக்குப் பொதுவான தனிச்சொற்கள். ஆள்களங்களின் எண்ணிக்கை அதிகமாக, அதிகமாக பொதுச் சொற்கள் குறைந்து கொண்டே வருகின்றன.

ஆள்களம்	தனிச் சொற்கள்	முழுவதற்கும் பொதுவான சொற்கள்
சமயம்	169277	169277
திரைப்படம்	103163	11587
கலை	222333	9960
விளையாட்டு	82384	5887
கவிதை	17065	1854
அறிவியல்	136931	1705
கல்வி	23749	1046
சரித்திரம்	241082	1044
புவியியல்	17051	711
சட்டம்	45738	680
மற்ற பிற	166283	680
மேலாண்மை	23116	589
ஜோதிடம்	13194	503
வேளாண்மை	12396	421
செய்தி	408833	420
பொருளாதாரம்	17873	392
அறிவியல்	227952	392
மொழியியல்	614184	392
அரசியல்	285040	392

என்பது தமிழ் மொழிக்குச் சரிவராது. எங்கள் ஆய்வில், எமில்லெ (Emille) பனுவலிலிருந்து தேர்ந்தெடுக்கப்பட்ட 2.4 லட்சம் தனிச்சொற்கள் கொண்ட பேரகராதியில் நாங்கள் கையெழுத்து உணர்வியைச் சோதனை செய்ய உபயோகித்த சாதாரணமான 2000 வார்த்தைகளில் 1530 வார்த்தைகளே இருந்தன [5] என்பது குறிப்பிடக்கூடு. ஈரெழுத்துப் (bigram) படிமங்களால் மட்டுமே கையெழுத்து உணர்வியின் (online handwriting recognition system) பிழைகளைத் திருத்த இயலும் [6]. அப்படியும் பல ஜோடி வார்த்தைகள் (உ: அவன், அவள்; வருவான், வருவாள்) கடைசி எழுத்தில் மட்டும் வேறுபடும் என்பதால், அவற்றிற்கு ஈரெழுத்துப் படிமங்களும் தவறான சொற்களைக் கொடுக்க வாய்ப்பு உள்ளது. ஆக, அத்தகைய சூழ்நிலையில் நாம் உணர்வியின் எழுதப்பட்ட சொல்லிலிருந்து கீற்றக்குறியீடுகளுக்கு (symbols) துண்டுபடுத்தல் (segmentation) [7] மற்றும் உணர்தல் (recognition) ஆகிய பணிகளின் துல்லியத்தையே முயன்று அதிகரிக்க வேண்டியிருக்கிறது [8]. இந்தக் கருத்து ஒளிவழி எழுத்துரு உணர்தல் (OCR) முறைமைகளுக்கும் [9, 10, 11] பொருந்தும்.



படம் 3. தமிழ் மொழியின் செறிவு : பனுவலில் ஒவ்வொரு ஆள்களமாகச் சேர்க்கும்போது (நீல நிறம்) தனிச் சொற்களின் எண்ணிக்கை தங்கு தடையின்றி வளர்தல். ஒவ்வொரு தடவையும் மற்றுமொரு லட்சம் பனுவல் சேர்க்கப்படும் பொழுது தனிச்சொற்களின் வளர்ச்சி (சிவப்பு நிறம்)

தமிழில் நீளமான சொற்கள் அதிகம் (Mean length of Tamil words is high)

தமிழில் ஒவ்வொரு வினைச்சொல்லுக்கும் குறைந்த பட்சமாக என்பது உருபனியல் அமைப்புகள் (paradigms) உள்ளன. கூட்டு வினைச்சொற்களையும் கணக்கில் கொண்டால் இது இன்னும் பல மடங்காக உயரும் [1]. தமிழ் உருபனியல் மாற்றங்கள் எழுவாய், பயனிலை, செயப்படுபொருளைக் கண்டறிவதன் மூலம் பிரதிபெயர்களின் ஆதியைக் கண்டறிதலில் (pronominal resolution) பயன்படுமளவுக்கு பொருள் செறிந்தது [12]. இத்தகைய செறிவு தெலுங்கு போன்ற மற்ற திராவிட மொழிகளிலும் உண்டு [13]. இச்செறிவே சொற்திருத்தி போன்ற முறைமைகளை உருவாக்குவதைக் கடினமாக்குகின்றன [14]. இவ்வித உருபனியல் செறிவானது ஒரு ஆங்கிலச் சொற்றொடரை மொழிபெயர்க்கையில் ஒரு வார்த்தையில் பொருள் தரும் வகையில் அமைந்துள்ளது [15]. உருபனியல் மாற்றச் செறிவினால் தமிழ் வார்த்தைகளின் சராசரி நீளம் அதிகமாக இருக்க வாய்ப்பு உள்ளது. இதனைச் சோதிப்பதற்காக வெவ்வேறு நீளமுள்ள தமிழ் தனி வார்த்தைகளின் புள்ளி விவரங்களைச் சேகரித்தோம். இவை அட்டவணை 4-ல் கொடுக்கப்பட்டுள்ளன. இங்கு நீளம் என்பது தமிழ்ச் சொல்லிலுள்ள உயிர்மெய் எழுத்துக்களின் எண்ணிக்கையே. உதாரணமாக, "கல்கி" என்ற சொல்லில் மூன்று எழுத்துக்கள் உள்ளன. ஆயின், ஒருங்குறி எண்ணிக்கை ஐந்தாகும் .

அட்டவணை 4: வெவ்வேறு நீளமுள்ள தமிழ் தனி வார்த்தைகளின் புள்ளி விவரங்கள் . மொத்த எண்ணிக்கையை ஒட்டி, இறங்கு வரிசைப்படி கொடுக்கப்பட்டுள்ளன .

எழுத்துக்கள் /சொல்	மொத்த எண்ணிக்கை	எழுத்துக்கள் /சொல்	மொத்த எண்ணிக்கை
10	135486	5	41230
11	130425	18	30471
9	128475	4	23679
12	122741	19	21826
8	115506	20	14660
13	107806	21	10195
7	90559	3	8465
14	90164	22	6818
15	72450	23	4512
6	64978	24	3035
16	55613	25	1924
17	41860	2	1833

முடிவுரை மற்றும் வருங்கால ஆய்வுப் பணிகள் (Conclusion and future work)

பை-கிராம் இடுகுறி அல்லது எழுத்து சார்ந்த புள்ளி விவரங்கள் வெவ்வேறு ஆள்களங்களுக்குக் குறிப்பிடத்தக்க அளவில் வேறுபடுகிறதா என்பதை உற்றறிதல் மூலம் பல பயனுறு தகவல்களைப் பெறலாம். உதாரணமாக, இப்புள்ளிவிவரங்கள் சங்கத்தமிழ்ப் பாடல்களுக்கும் நவீன தமிழ் தரவுகளுக்கும் குறிப்பிடத்தக்க விதத்தில் வேறுபட்ட தனிச்சிறப்பியல்புகளைக் கொண்டிருக்கும். மேலும் இலக்கண, சொற்பொருளியல் மற்றும் புணர்ச்சி விதிகளை வரையறுப்பதன் மூலம் உணர்தல் நெறிமுறைகளின் செயலாக்கம் வலுவடையலாம்.

நன்றி (Acknowledgment)

இந்தக் கட்டுரையிலுள்ள கருத்துக்கள், இந்திய அரசின் இந்திய மொழிகளுக்கான தொழில் நுட்ப வளர்ச்சித் திட்டத்தின் (Technology Development for Indian Languages – TDIL) நிதி உதவியுடன் நாங்கள் மேற்கொண்ட ஆய்வுச் செயல் திட்டங்களில் பணி புரியும்போது பெற்ற அனுபவங்கள், எதிர்கொண்ட இடையூறுகள் முதலியவற்றால் ஊக்குவிக்கப்பட்டன. அதற்காக, நாங்கள் இந்திய அரசின் தொலைத்தொடர்பு மற்றும் தகவல் தொழில் நுட்பத் துறைக்குக் கடமைப் பட்டுள்ளோம். மேலும், முதன்மையாக, எங்கள் திருக்குரல் ஒளிவழி எழுத்துரு உணர்தல் முறைமையைப் பயன்படுத்தும் பலரும் கொடுத்த தமிழ் உரைகளிலிருந்தே இந்தக் கட்டுரைக்கான பனுவல் தயாரிக்கப்பட்டது. அதற்காக, அவர்கள் அனைவருக்கும் எங்கள் நன்றி.

மற்றும் பேராசிரியர் தெய்வசுந்தரம் அவர்களும், முனைவர். வாசு ரெங்கநாதன் அவர்களும் இந்தக்கட்டுரையை நாங்கள் வழங்கும் முறையின் தரத்தை உயர்த்த ஆலோசனைகள் வழங்கினார்கள். அவர்களுக்கும் எங்கள் நன்றிகள் பல.

மேற்கோள்கள் (References)

1. Rajendran, S., Viswanathan, S. Ramesh Kumar, "Computational morphology of Tamil verbal complex," Language in India, Vol. 3 : 4 April 2003.
2. Annamalai, E. Dynamics of Verbal Extension in Tamil. Thiruvananthapuram: Dravidian Linguistics Association, 1985.
3. Richard Lederer. A celebration of English, good grammar, and wordplay. Marion Street Press, 2012.
4. David Crystal, "The Stories of English", Overlook TP Publishers, 2004. ISBN:978-1-58567-719-1.
5. Suresh Sundaram, Bhargava Urala and A. G. Ramakrishnan, "Language models for online handwritten Tamil word recognition," Proc. Workshop on Document Analysis and Recognition (DAR 2012), 16 December 2012, IIT Bombay, Mumbai, India.
6. Suresh Sundaram and A. G. Ramakrishnan, "Bigram language models and reevaluation strategy for improved recognition of online handwritten Tamil words," revised manuscript under review, ACM Transactions on Asian Language Information Processing (TALIP), 2013.
7. Suresh Sundaram and A. G. Ramakrishnan, "Attention-feedback based robust segmentation of online handwritten isolated Tamil words," ACM Transactions on Asian Language Information Processing (TALIP), Vol. 12 (1), March 2013, Article No. 4.
8. Suresh Sundaram and A. G. Ramakrishnan, "Performance enhancement of online handwritten Tamil symbol recognition with reevaluation techniques," revised manuscript under review, Pattern Analysis and Applications, 2013.
9. K.G.Aparna and A.G.Ramakrishnan, "A complete Tamil Optical Character Recognition System," Proc. Fifth IAPR Workshop on Document Analysis Systems DAS-02, Princeton, NJ, August 19-21, 2002, pp. 53-57.
10. A.G.Ramakrishnan and Kaushik Mahata, "A Complete OCR for Printed Tamil Text," Proc. Tamil Internet 2000, Singapore, July 22-24, 2000, pp. 165-170.
11. D.Dhanya and A.G.Ramakrishnan, "Simultaneous recognition of Tamil and Roman scripts," Proc. Tamil Internet 2001, Kuala Lumpur, August 26-28, 2001, pp. 64-68.
12. Murthy, Kavi Narayana, L. Sobha, and B. Muthukumari, "Pronominal resolution in Tamil using machine learning," Proc. First Intern. Workshop on Anaphora Resolution (WAR-I), pp. 39-50. 2007.
13. K. Narayana Murthy, "Parsing Telugu in the UCSG Formalism," Proc. Indian Congress on Knowledge and Language, vol 2, Jan. 1996, pp 1-16, Central Institute of Indian Languages, Mysore.
14. Ranjani Parthasarathy and Geetha T.V., Morphological Analyzer for Tamil, TI 2001, Chennai.
15. Vu, Ngoc Thang, and Tanja Schultz, "Initial experiments with Tamil LVCSR," Proc. Intern. Conf. Asian Language Processing (IALP), Hanoi, Vietnam. pp. 81-84, IEEE, 2012.