

Inaudible Voice Commands: The Long-Range Attack and Defense



Nirupam Roy



Sheng Shen



Haitham Hassanieh



Romit Roy Choudhury

University of Illinois at Urbana-Champaign



50 million voice assistants are sold in US



Inaudible Acoustics



Normal Sound
(< 24 kHz)



Ultrasound
(> 25 kHz)



“Inaudible Acoustics”
(> 25 kHz)



“Alexa, open the garage door!”



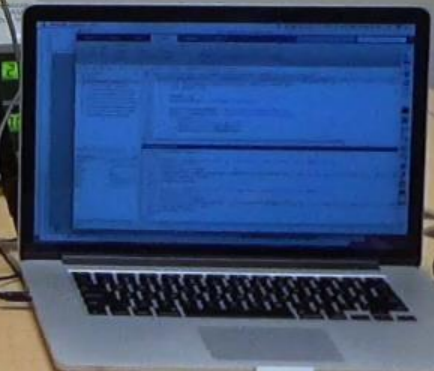
Inaudible Acoustics



Human speaker

Inaudible speaker

Alexa (Amazon Echo)



Talk Outline

0. [BackDoor], [DolphinAttack], [Princeton Video]

MobiSys'17
(Best Paper)

CCS'17

arXiv

Talk Outline

0. [BackDoor], [DolphinAttack], [Princeton Video]

Today's Talk:

1. How to launch long-range (realistic) attacks?

2. How to defend against these attacks?

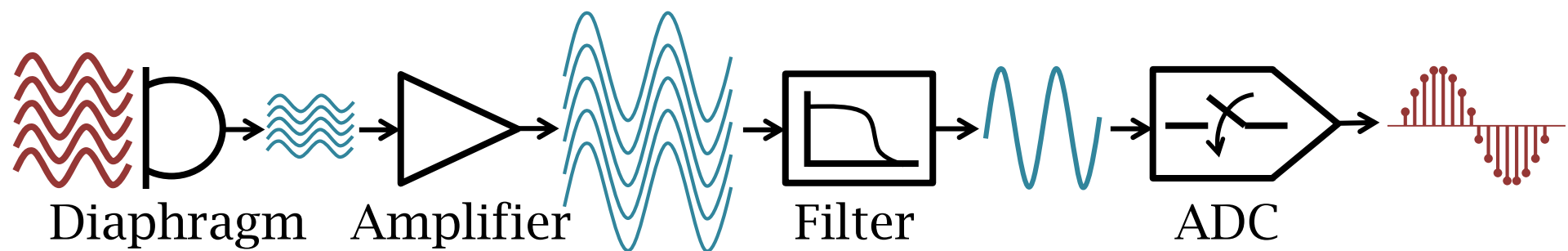
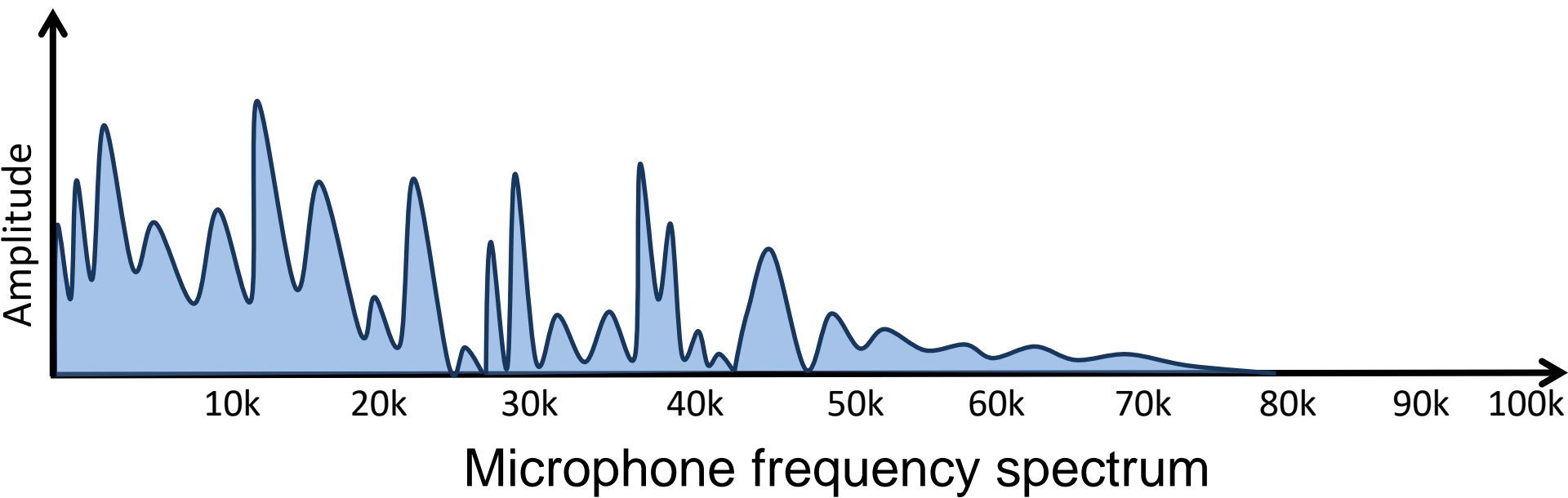
Talk Outline

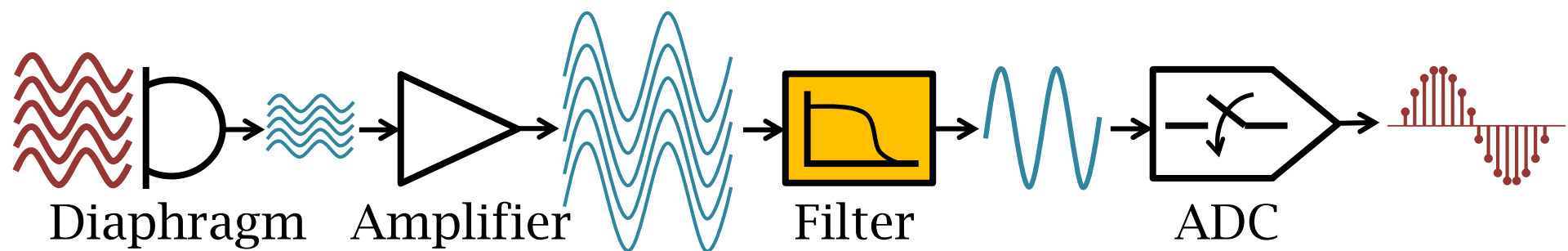
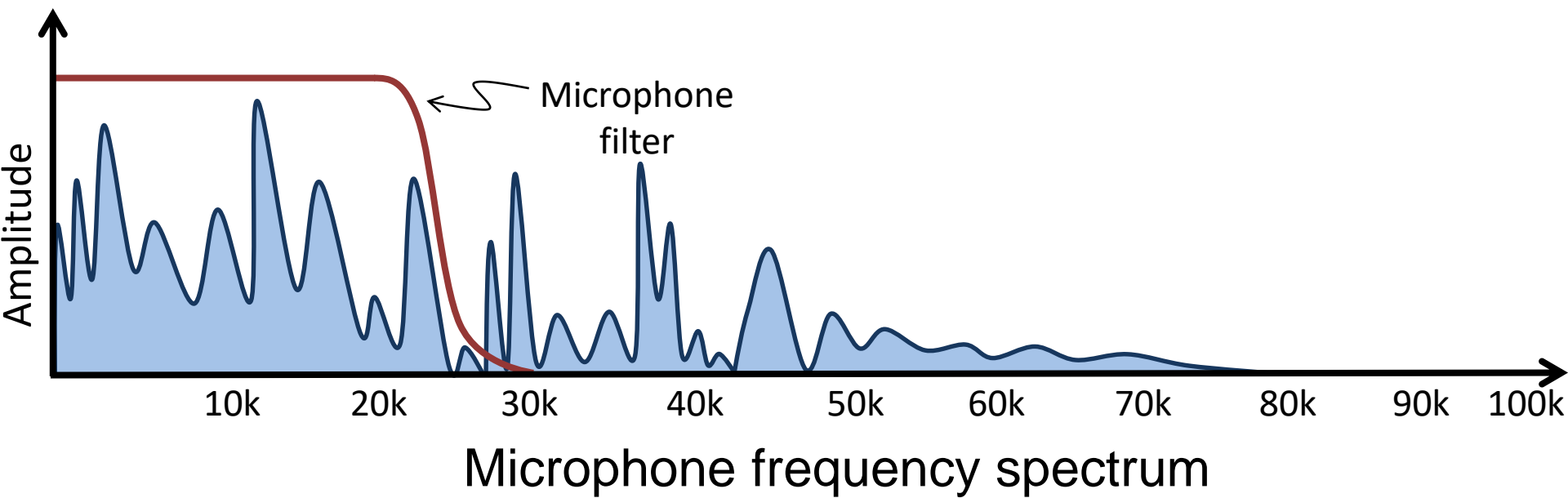
0. [BackDoor], [DolphinAttack], [Princeton Video]

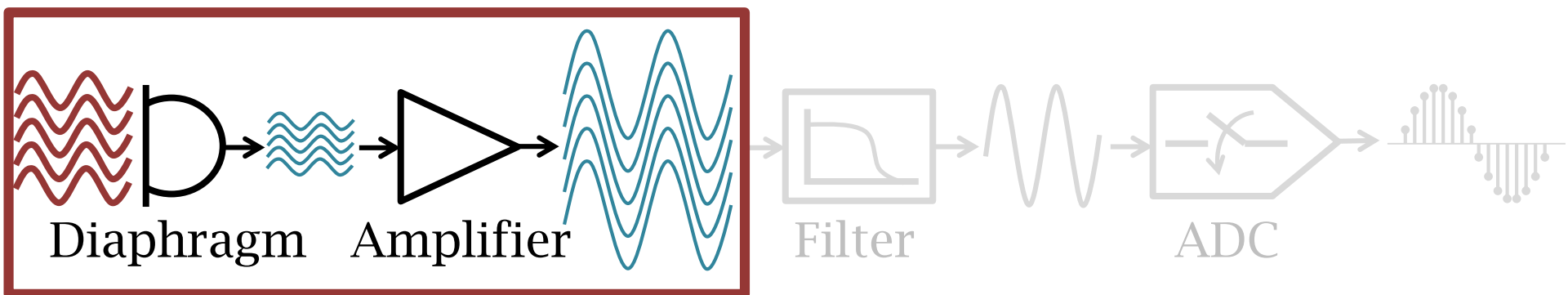
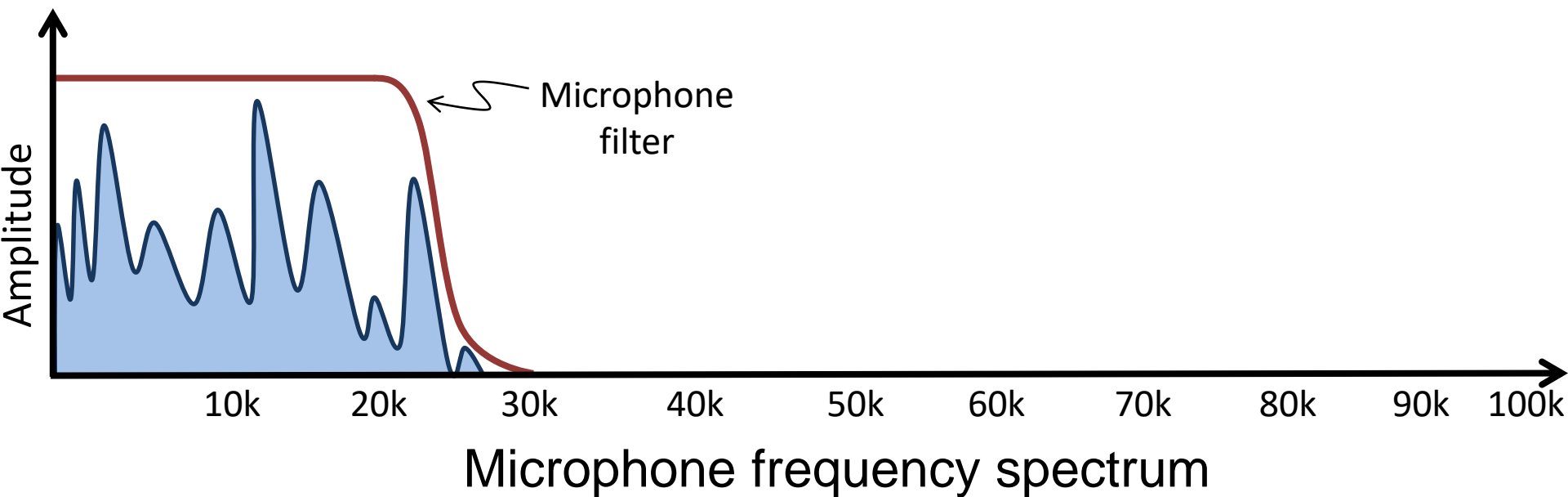
Today's Talk:

1. How to launch long-range (realistic) attacks?

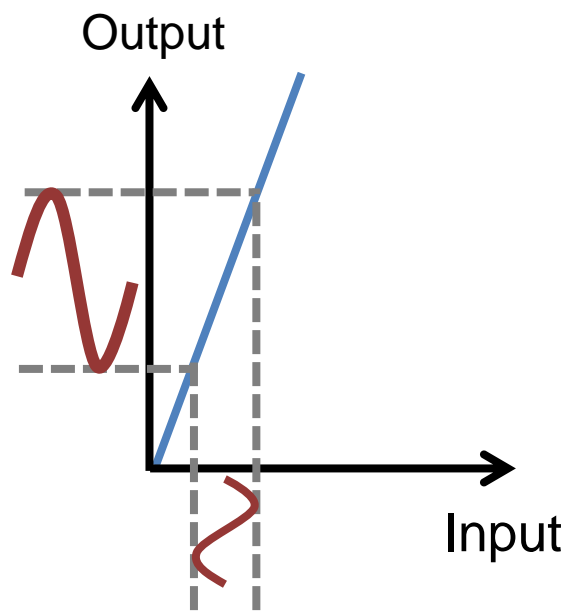
2. How to defend against these attacks?



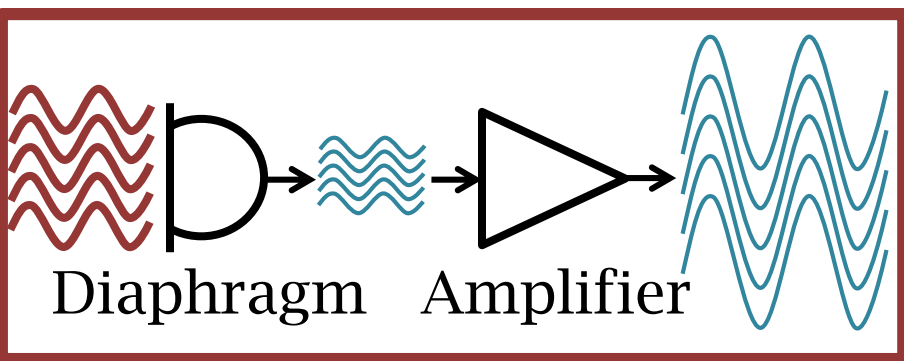


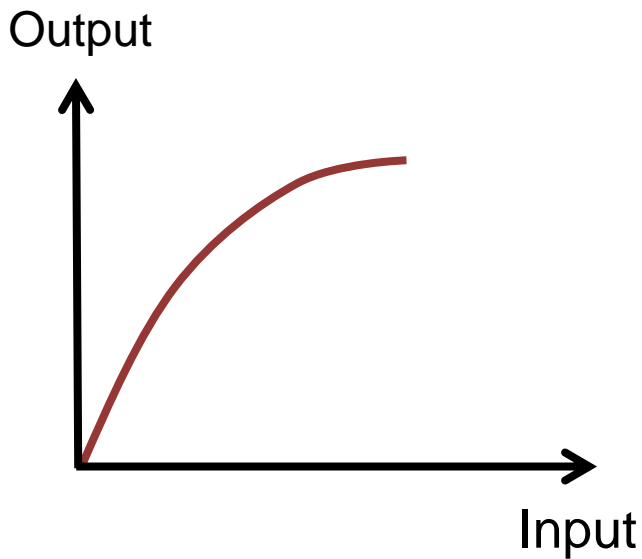
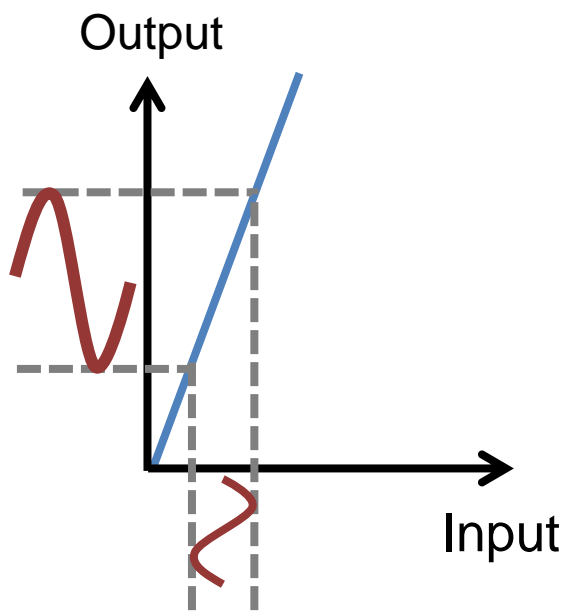


Air Vibration → Electric Voltage



$$V_{out} = a_1 V_{in}$$





Nonlinear

$$V_{out} = a_1 V_{in}$$

$$V_{out} = a_1 V_{in} + a_2 V_{in}^2 + \dots$$

Frequency

10k

20k

30k

40k

50k

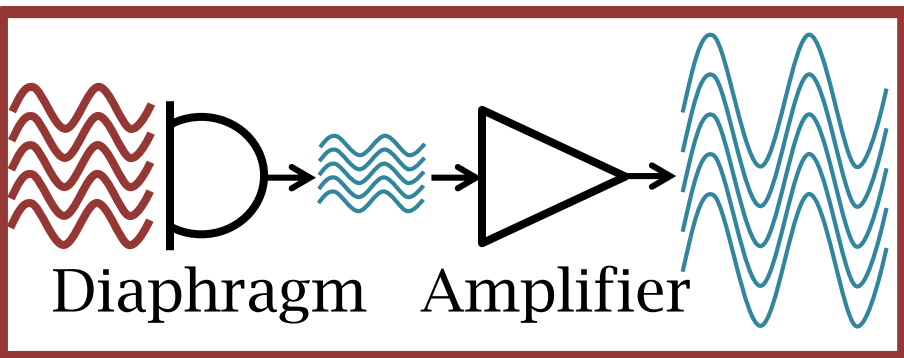
60k

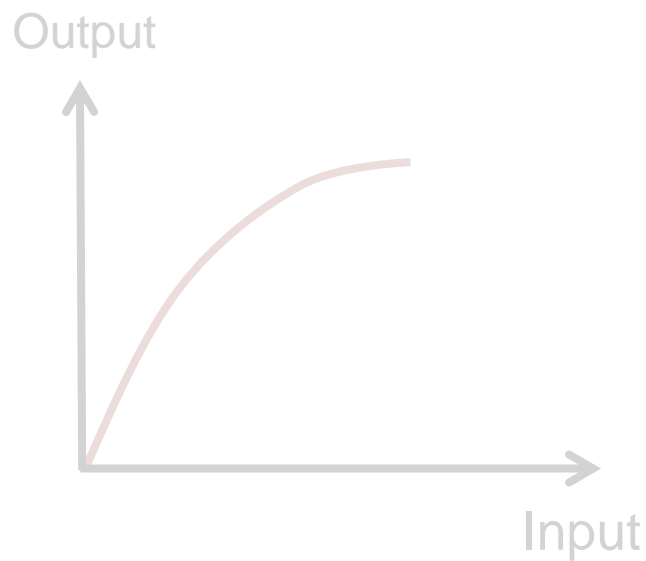
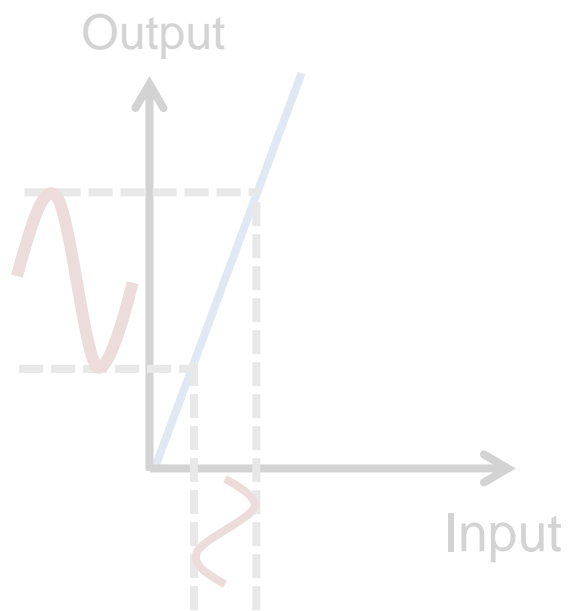
70k

80k

90k

100k





Nonlinear

$$V_{out} = a_1 V_{in}$$

$$V_{out} = a_1 V_{in} + a_2 V_{in}^2 + \dots$$

Frequency

10k

20k

30k

40k

50k

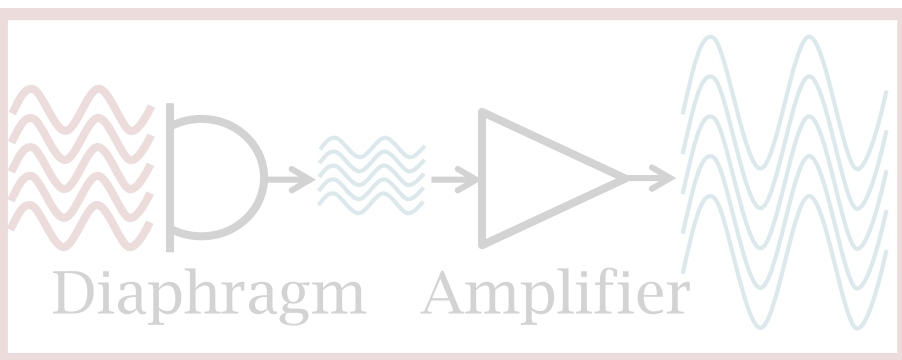
60k

70k

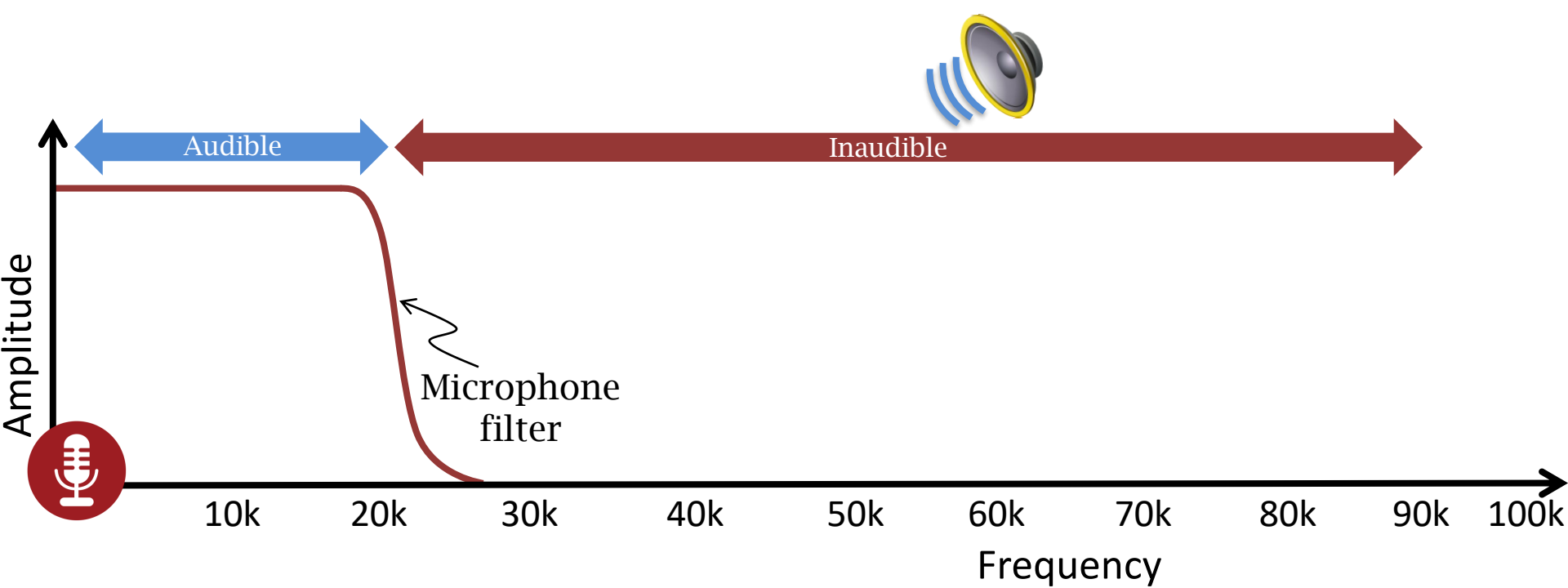
80k

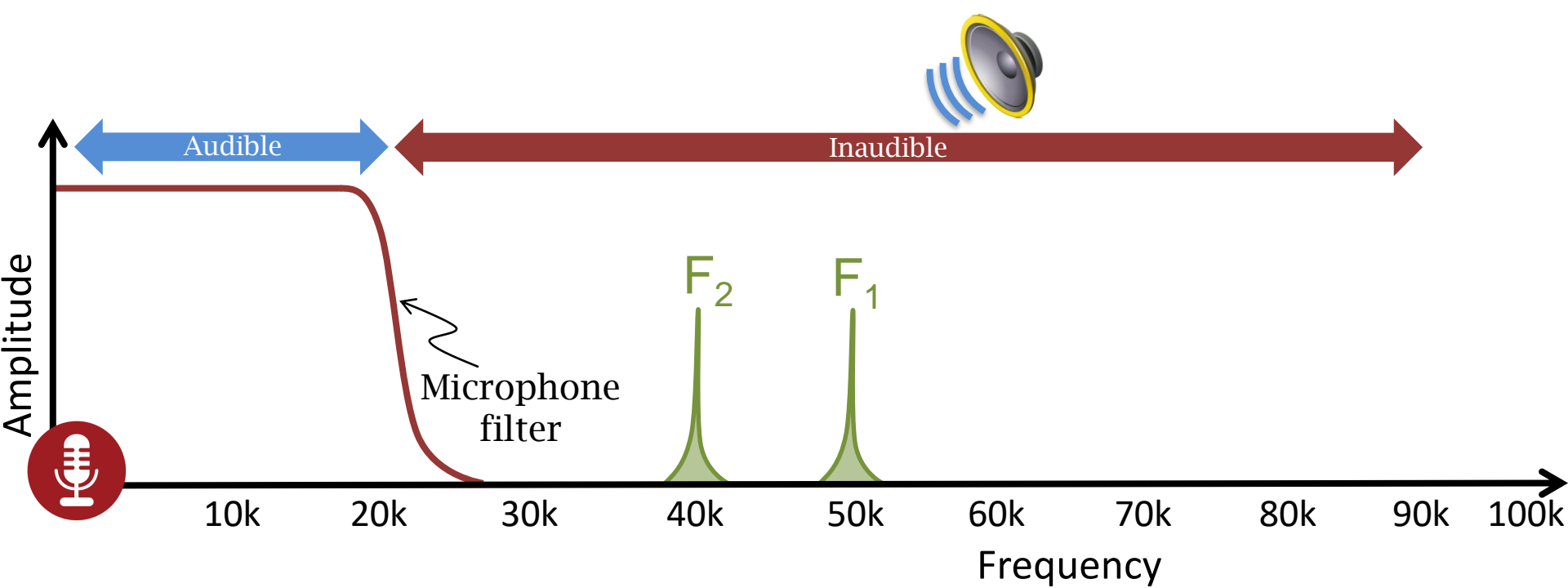
90k

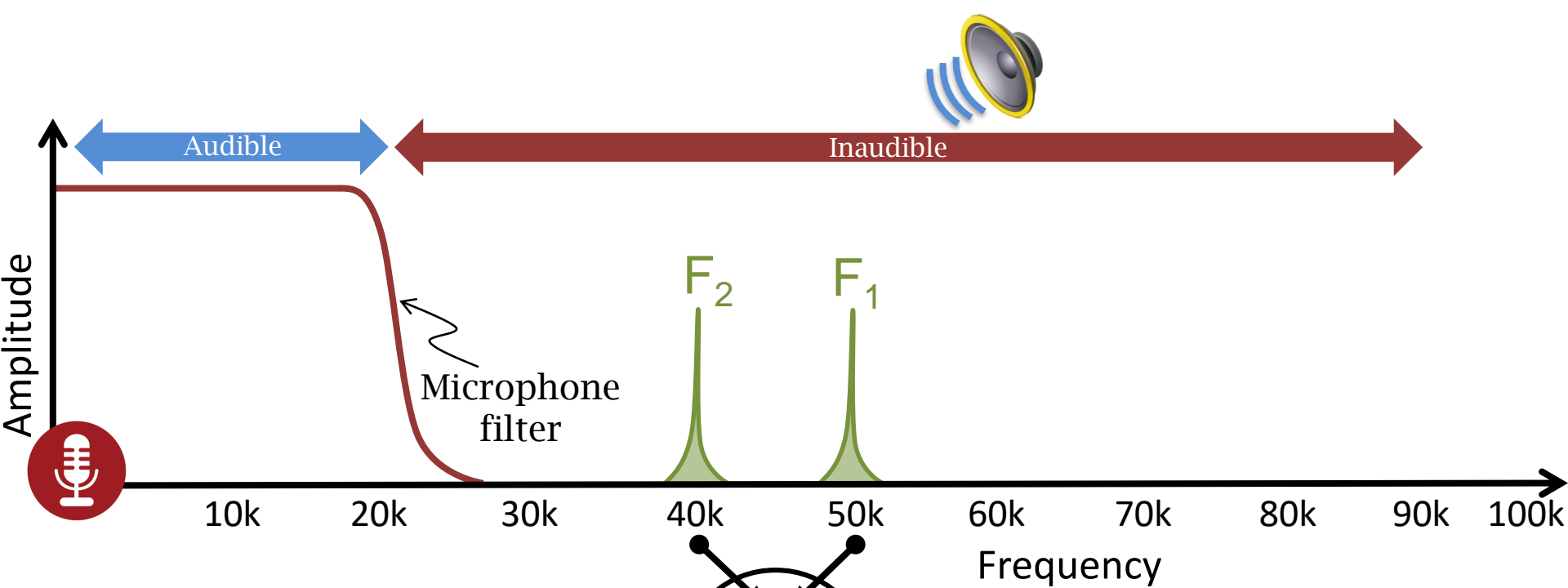
100k





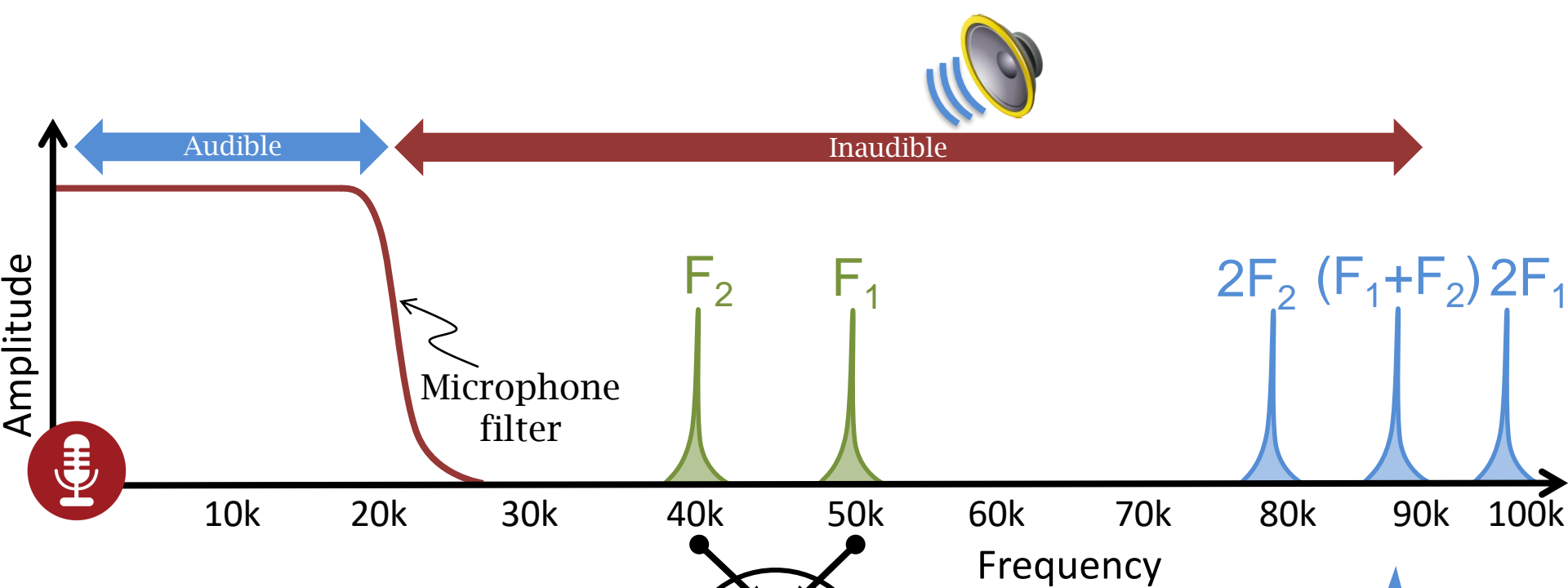






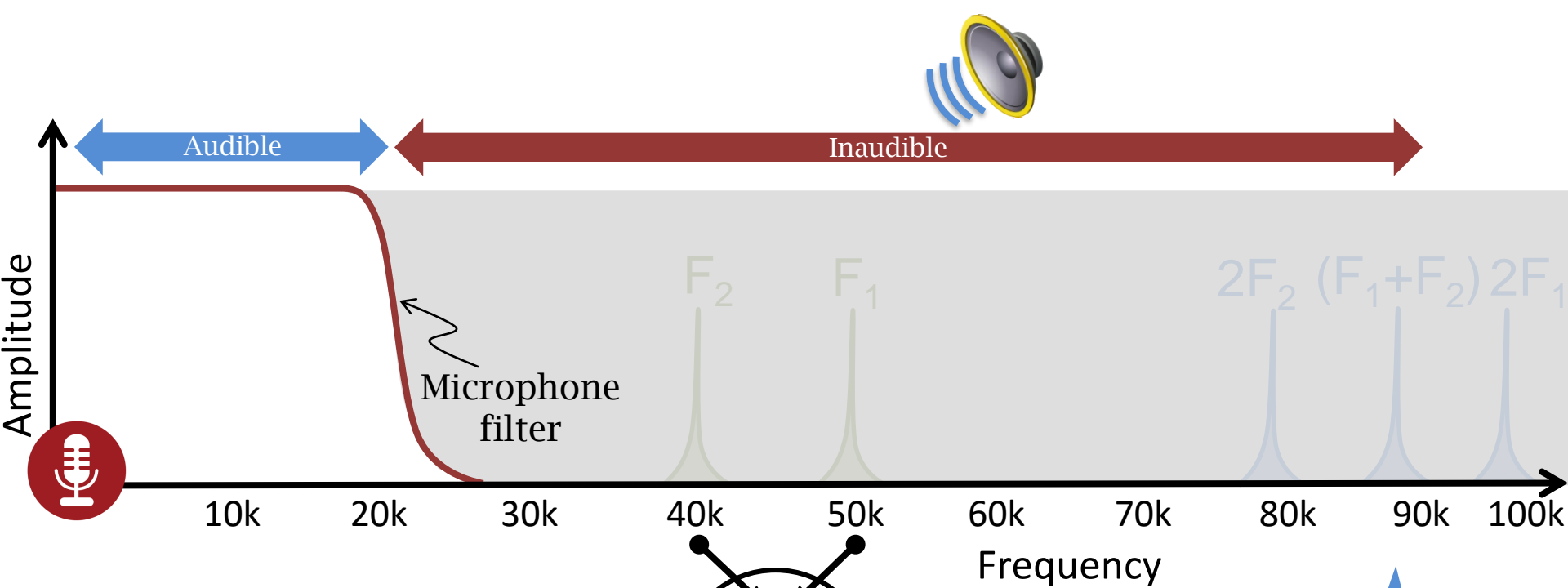
$$V_{out} = a_1 V_{in} + a_2 V_{in}^2$$

$$\begin{aligned} (\sin F_1 + \sin F_2)^2 = & -\cos 2F_1 \\ & -\cos 2F_2 \\ & -\cos (F_1 + F_2) \\ & +\cos (F_1 - F_2) \end{aligned}$$



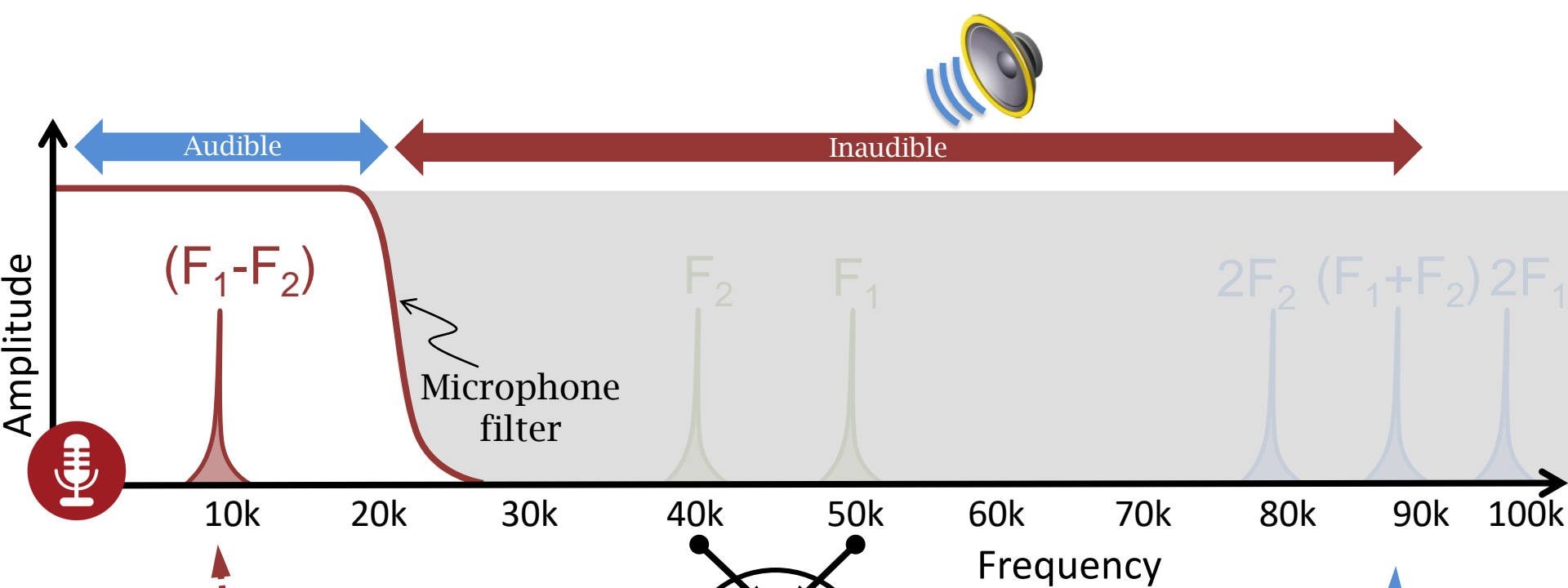
$$V_{out} = a_1 V_{in} + a_2 V_{in}^2$$

$$(\sin F_1 + \sin F_2)^2 = -\cos 2F_1 - \cos 2F_2 - \cos (F_1 + F_2) + \cos (F_1 - F_2)$$



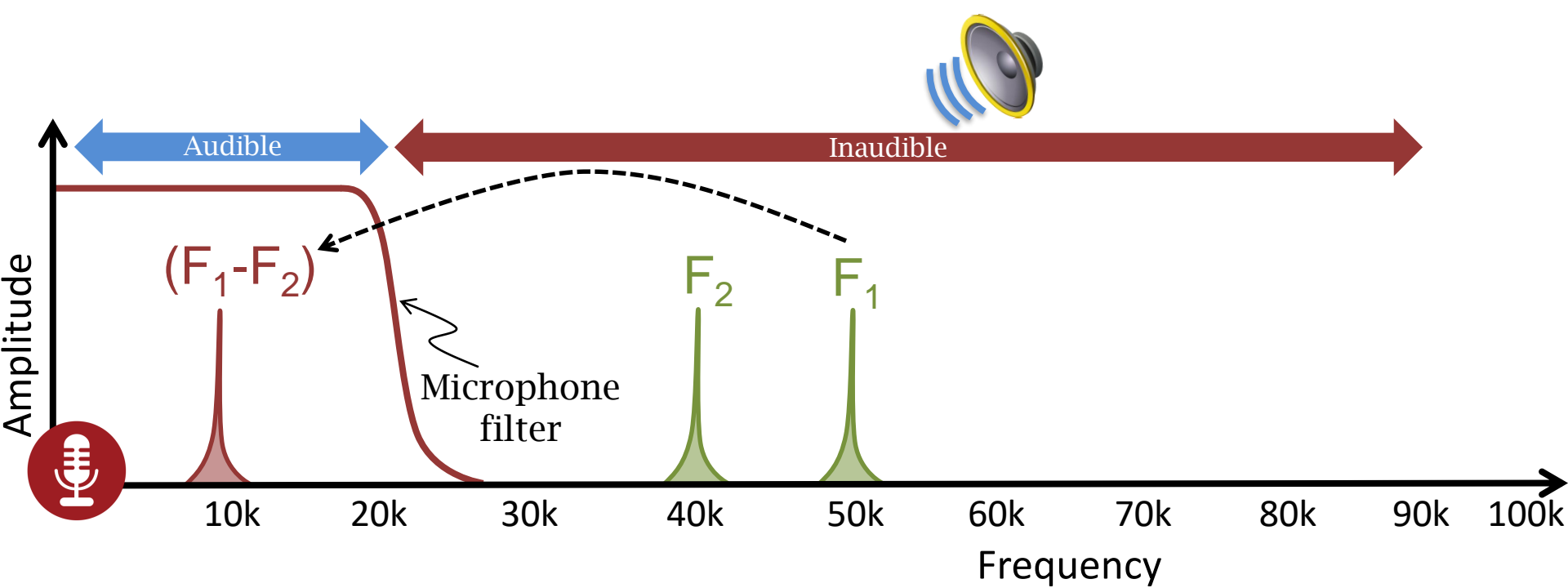
$$V_{out} = a_1 V_{in} + a_2 V_{in}^2$$

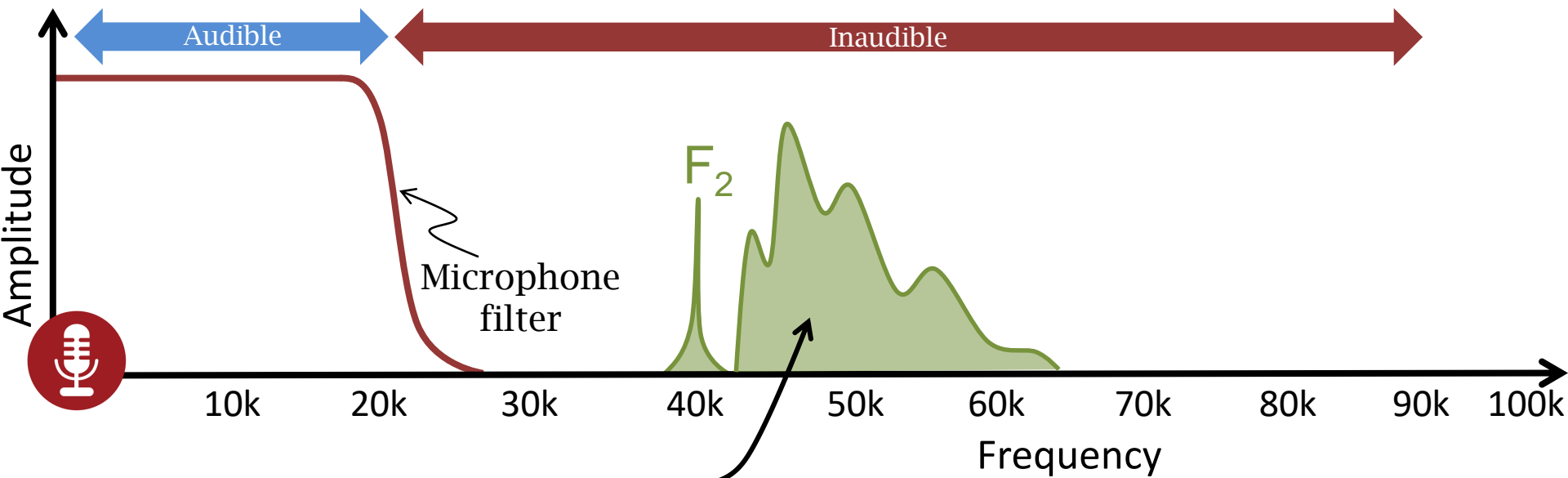
$$(\sin F_1 + \sin F_2)^2 = -\cos 2F_1 - \cos 2F_2 - \cos (F_1 + F_2) + \cos (F_1 - F_2)$$



$$V_{out} = a_1 V_{in} + a_2 V_{in}^2$$

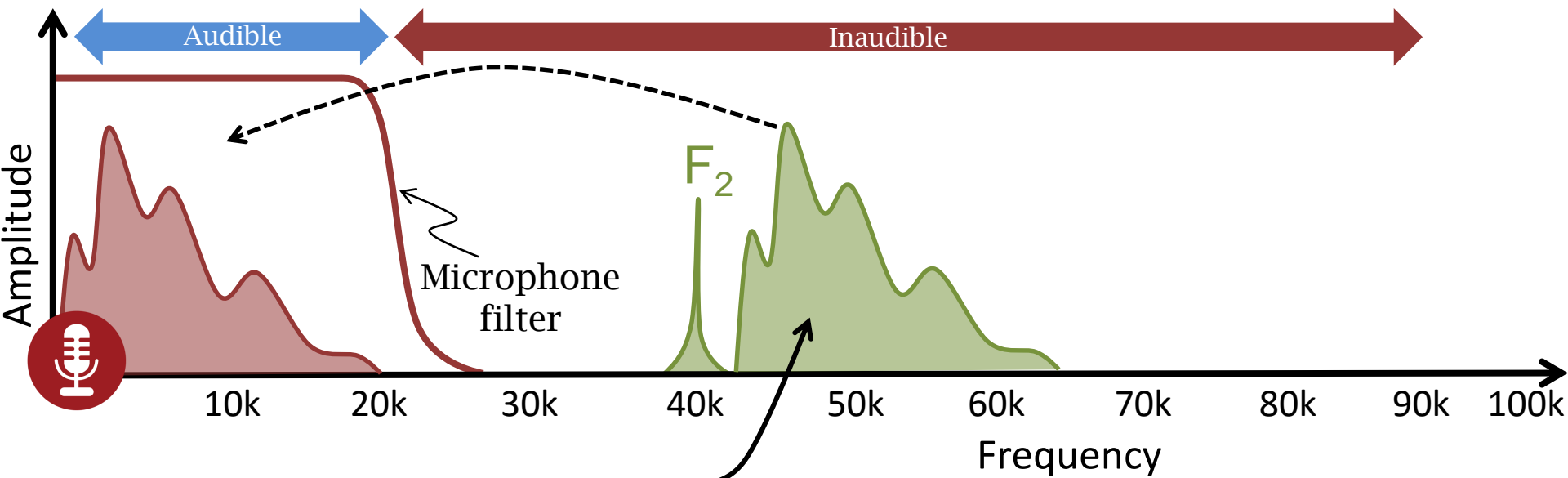
$$(\sin F_1 + \sin F_2)^2 = -\cos 2F_1 - \cos 2F_2 - \cos (F_1 + F_2) + \cos (F_1 - F_2)$$





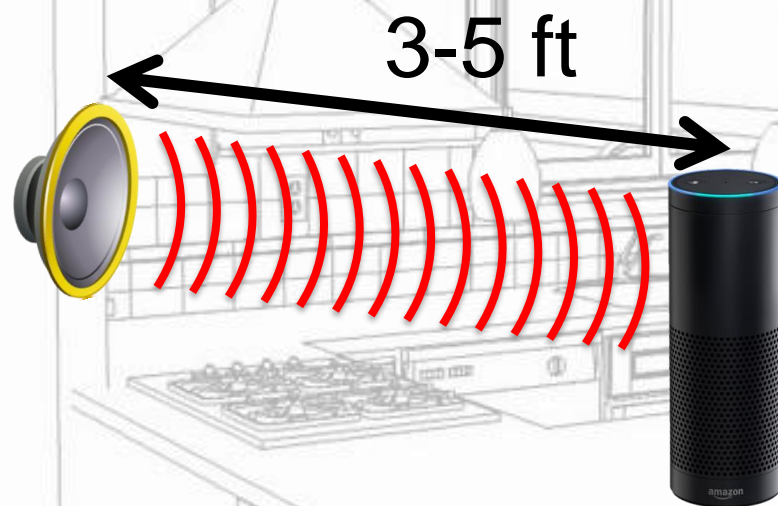
$V(t) = \text{"Alexa, open the garage door!"}$



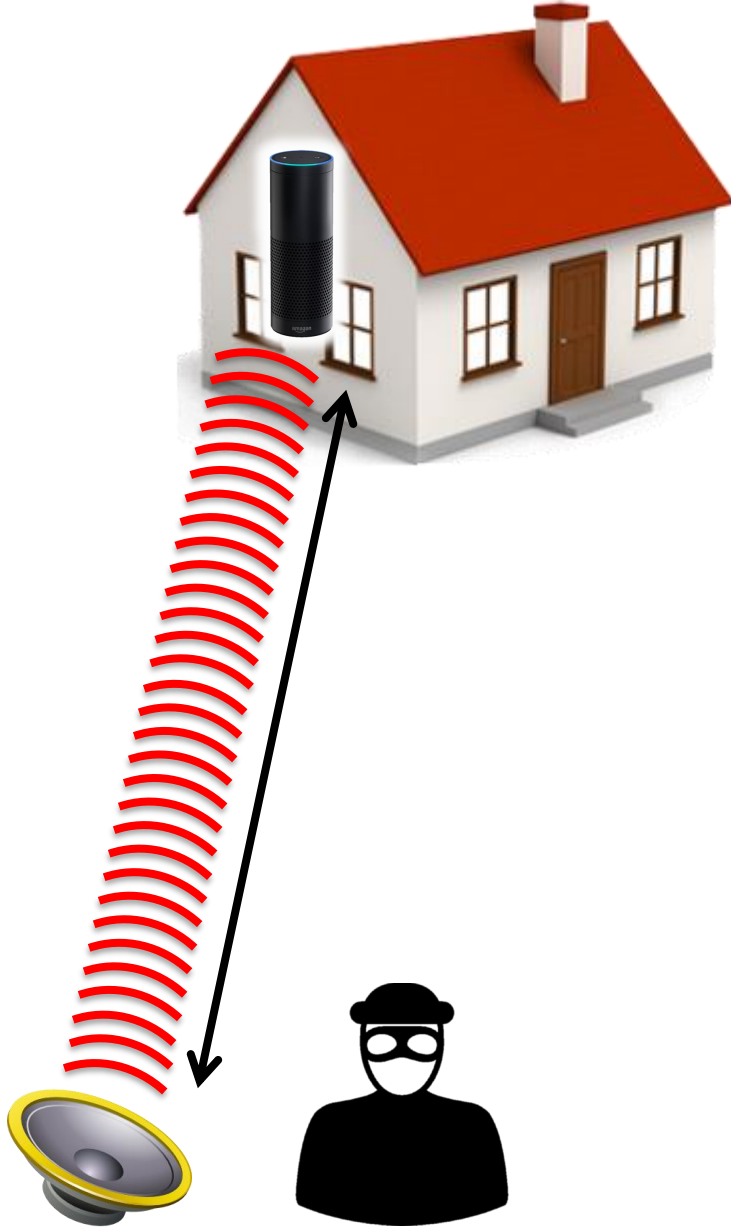


$V(t)$ = "Alexa, open the garage door!"

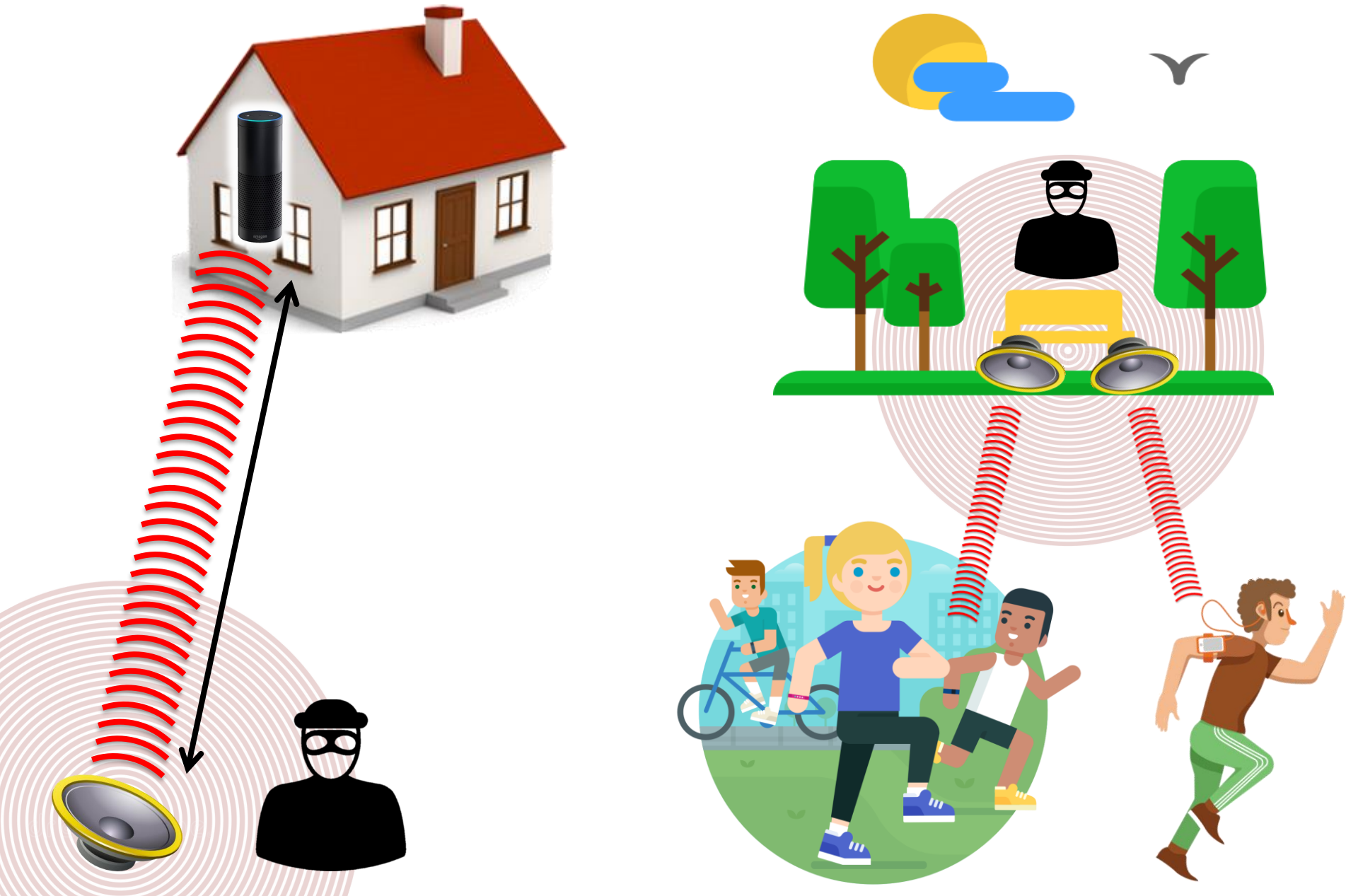




Can someone attack from a longer range?



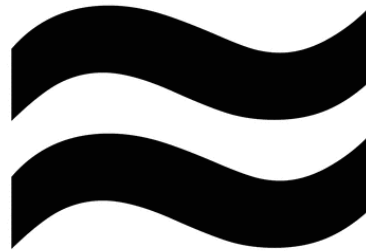
Can someone attack from a longer range?



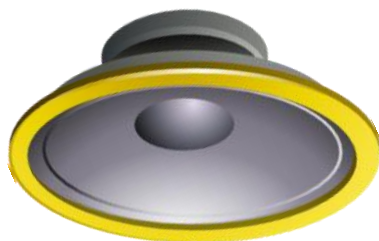
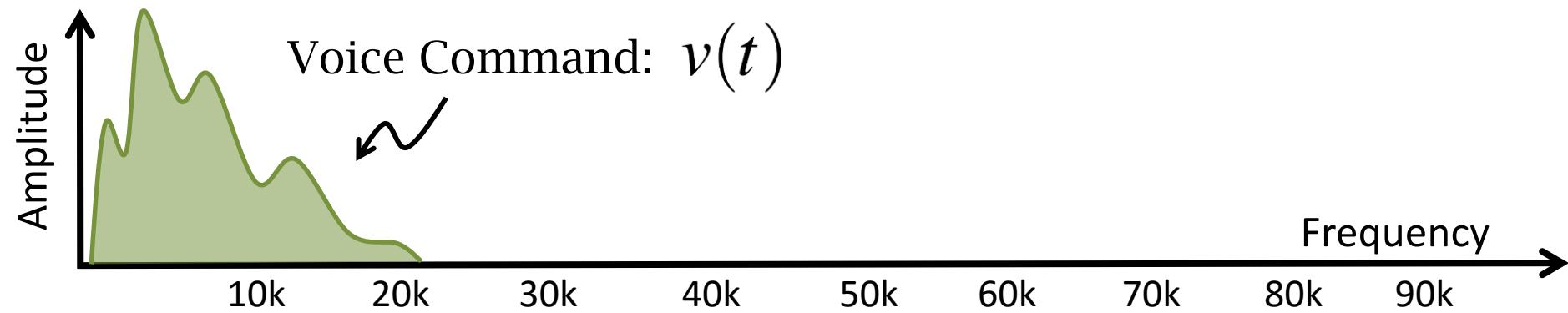
Can someone attack from a longer range?

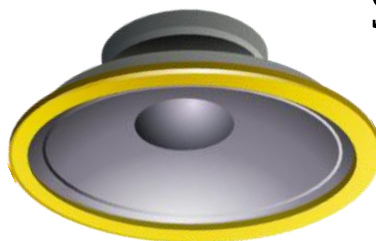
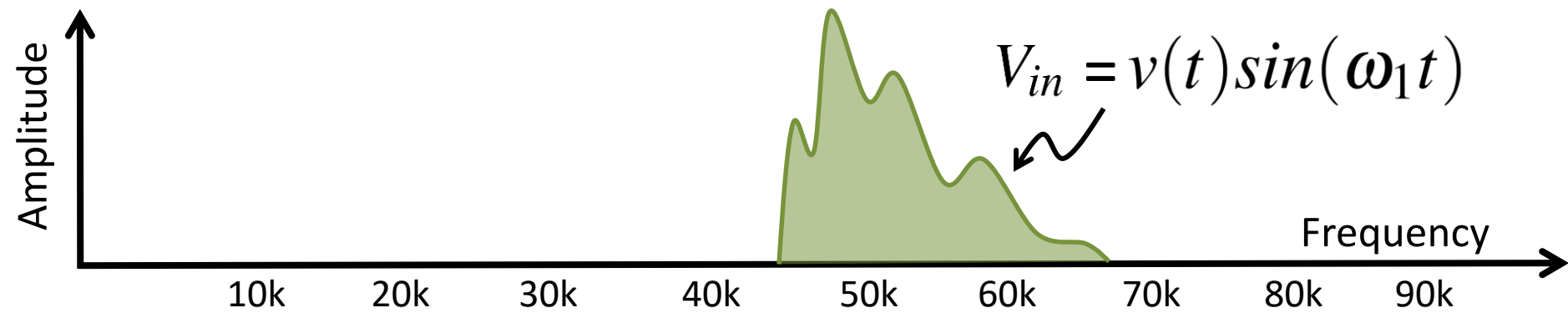


High power makes ultrasonic speakers audible



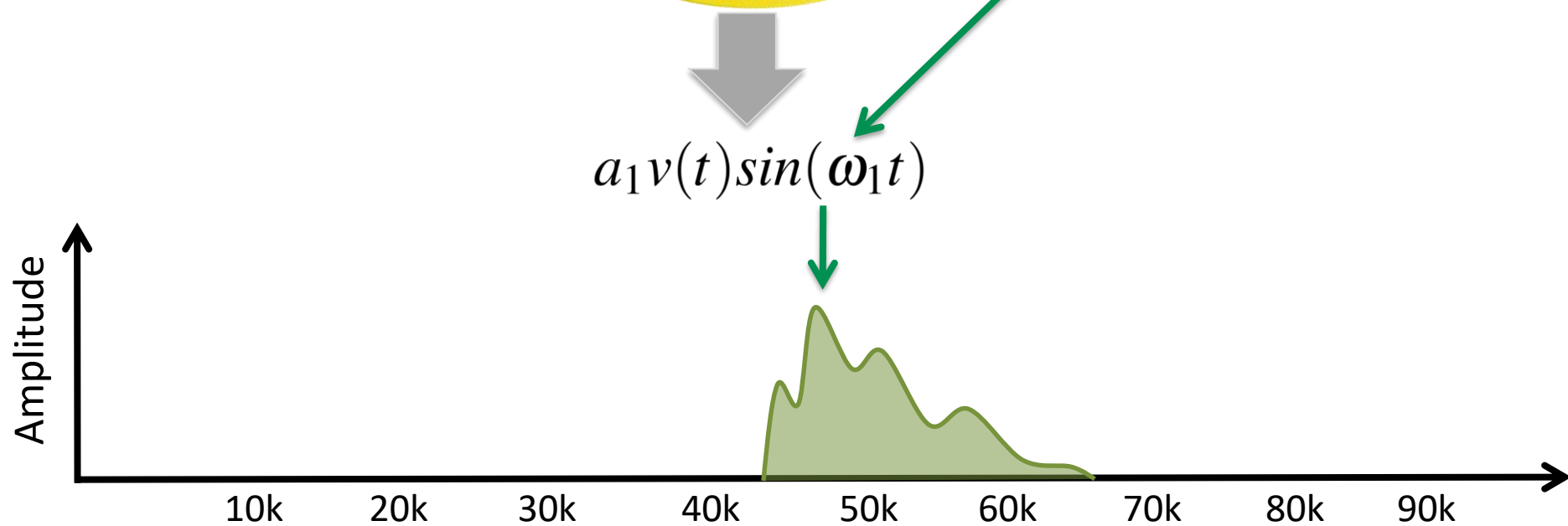
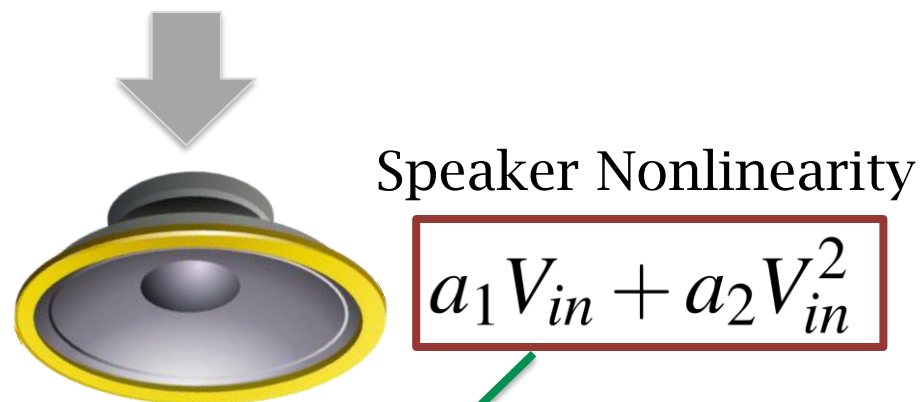
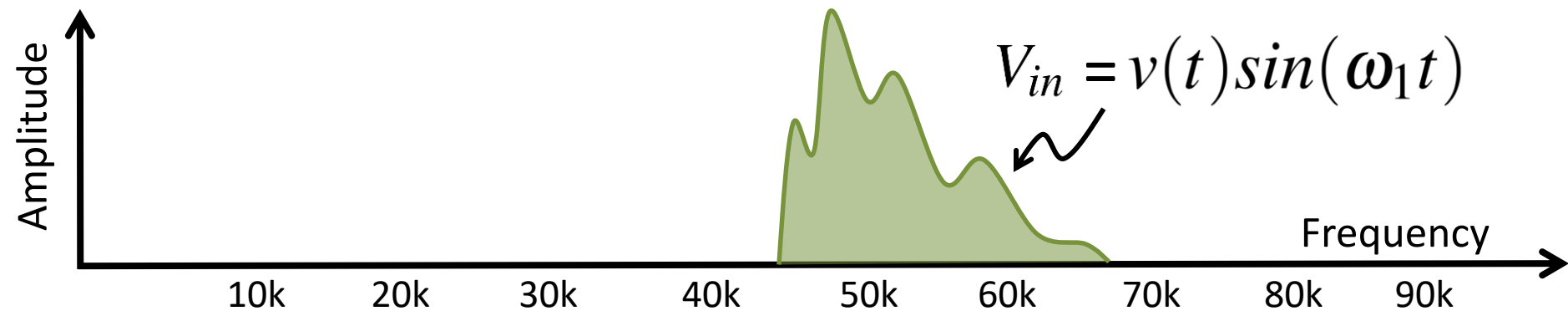
Speakers have nonlinearity too!

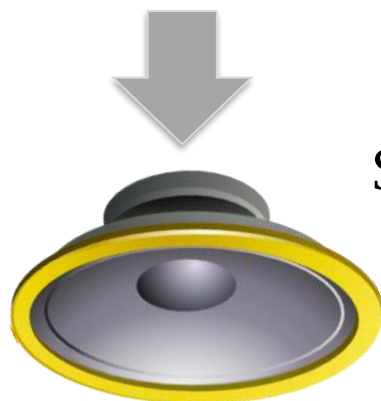
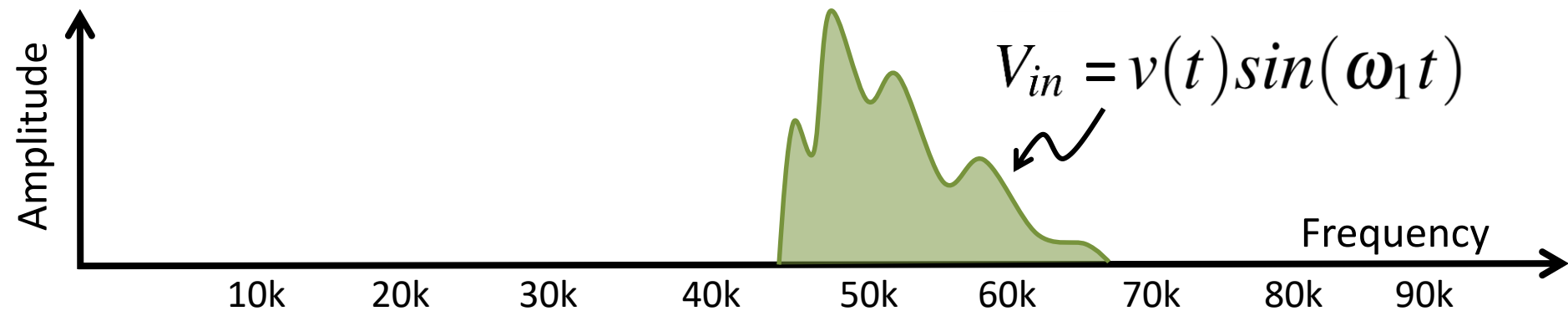




Speaker Nonlinearity

$$a_1 V_{in} + a_2 V_{in}^2$$

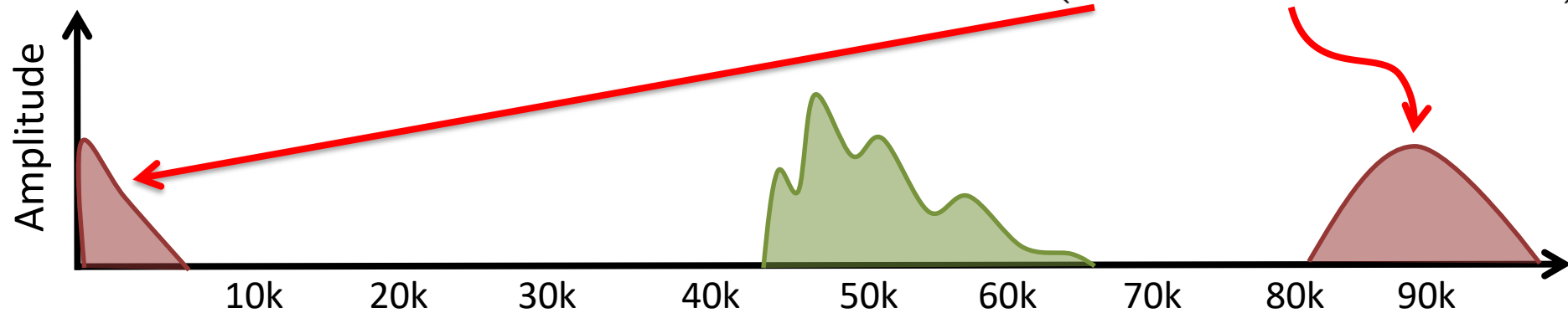


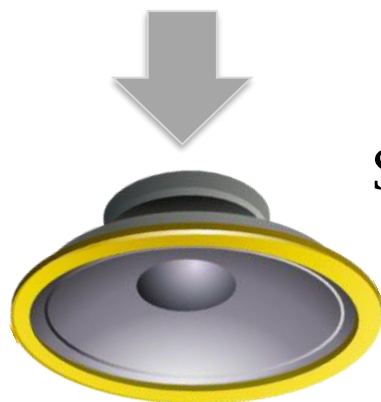
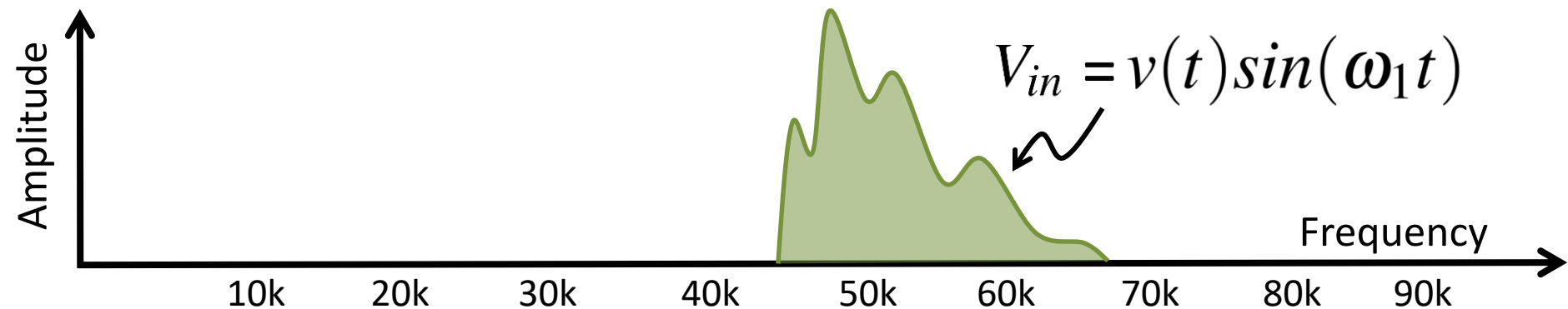


Speaker Nonlinearity

$$a_1 V_{in} + a_2 V_{in}^2$$

$$a_1 v(t)\sin(\omega_1 t) + a_2 \left(v^2(t) - v^2(t)\cos(2\omega_1 t) \right)$$





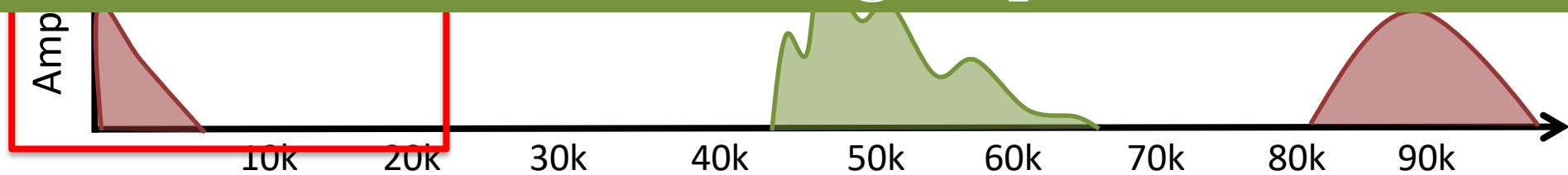
Speaker Nonlinearity

$$a_1 V_{in} + a_2 V_{in}^2$$

General to all speakers!

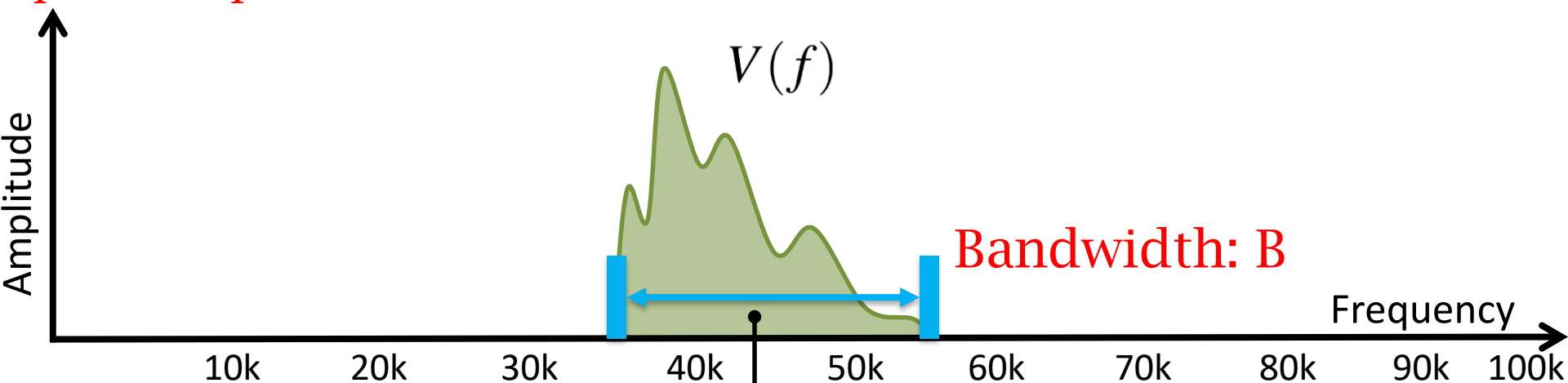


Our Solution: “Leakage Optimization”

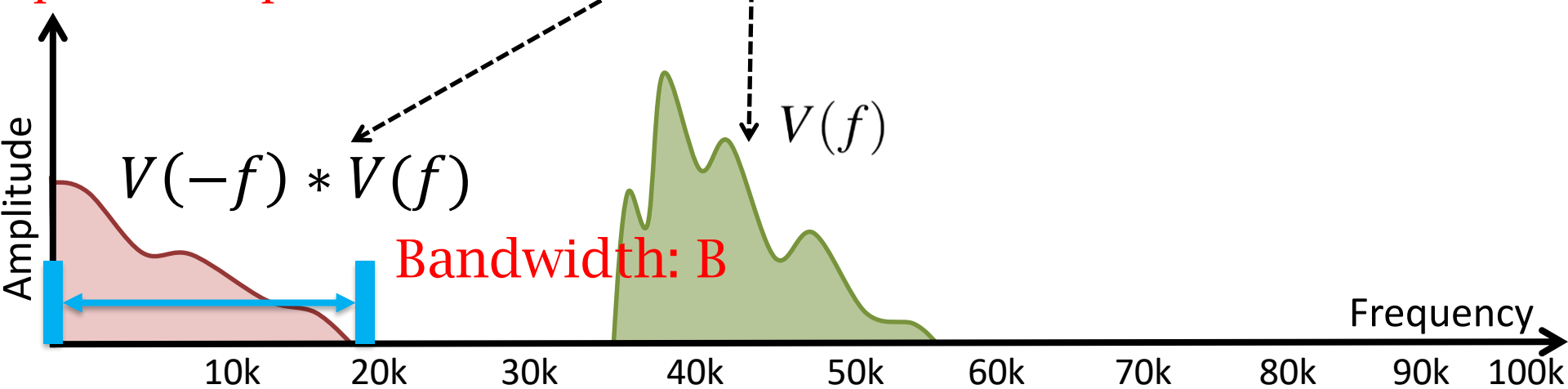


Speaker Nonlinearity \rightarrow Audible Leakage

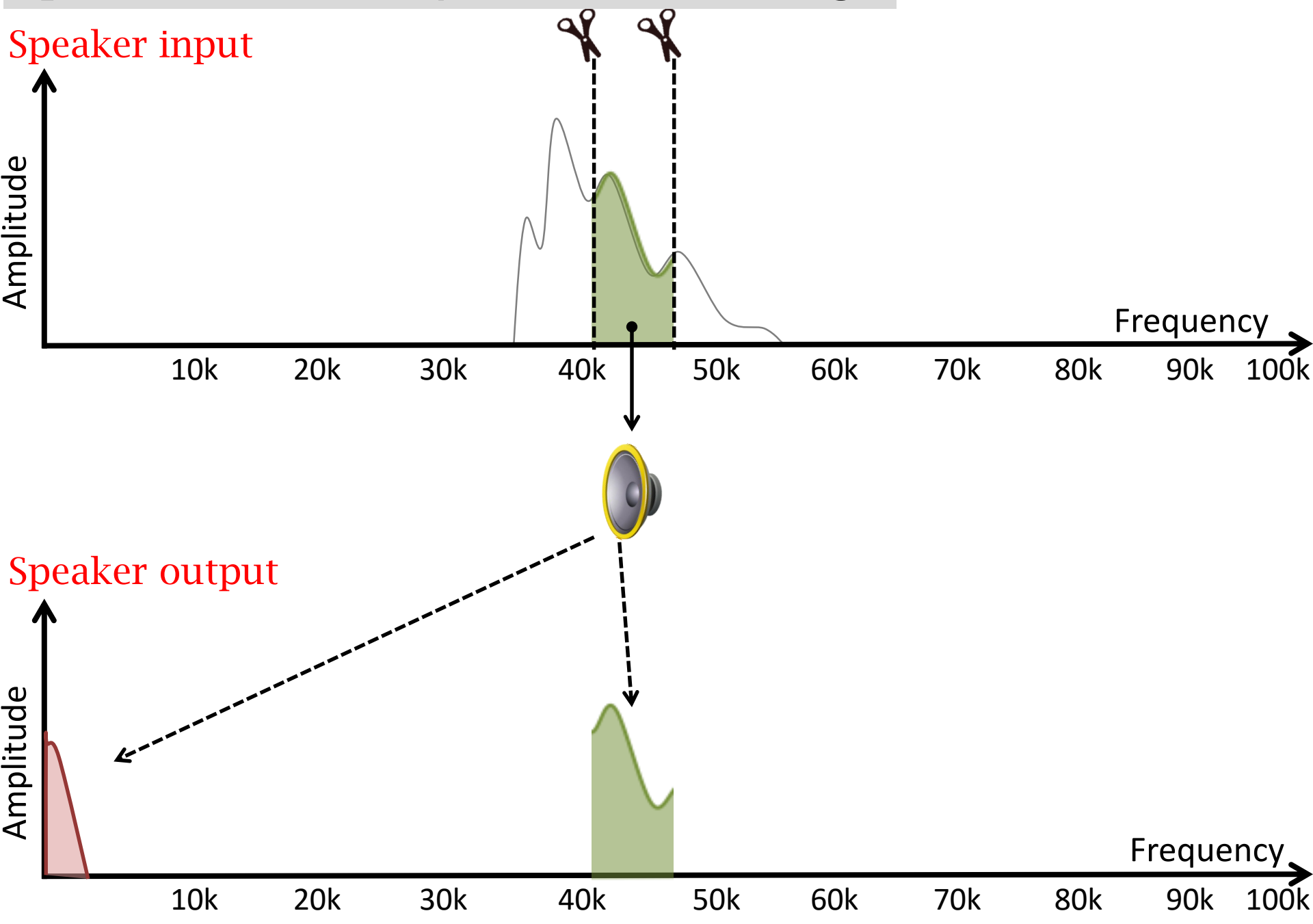
Speaker input



Speaker output

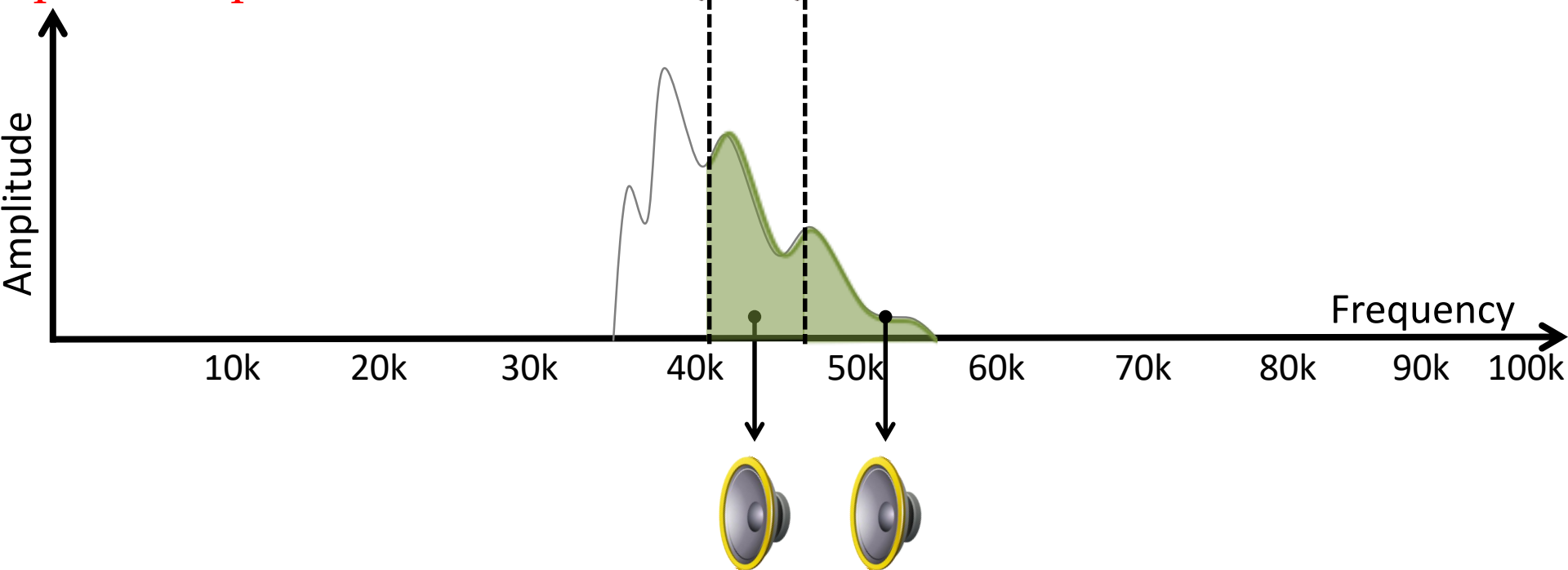


Speaker Nonlinearity \rightarrow Audible Leakage

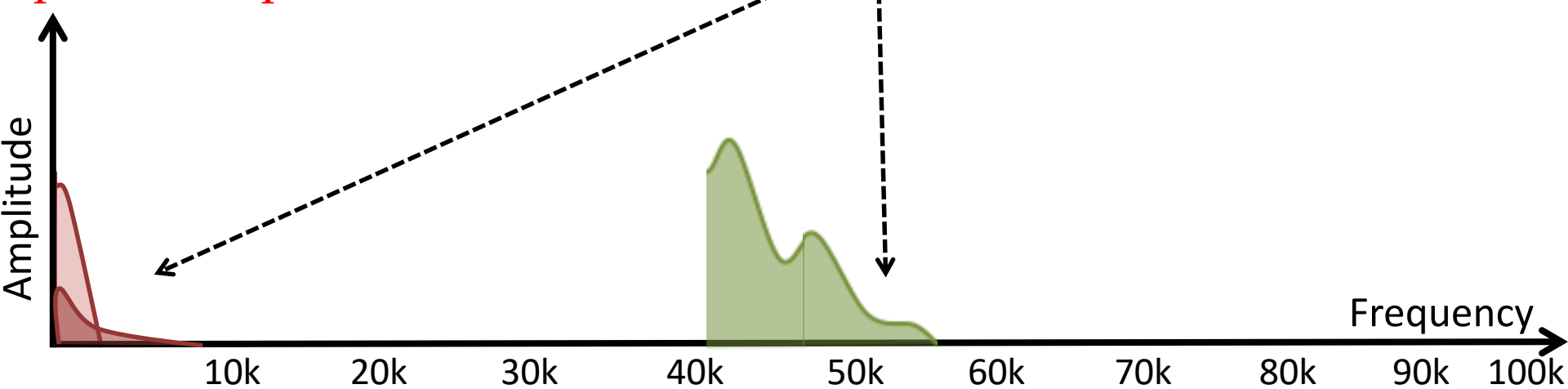


Speaker Nonlinearity \rightarrow Audible Leakage

Speaker input

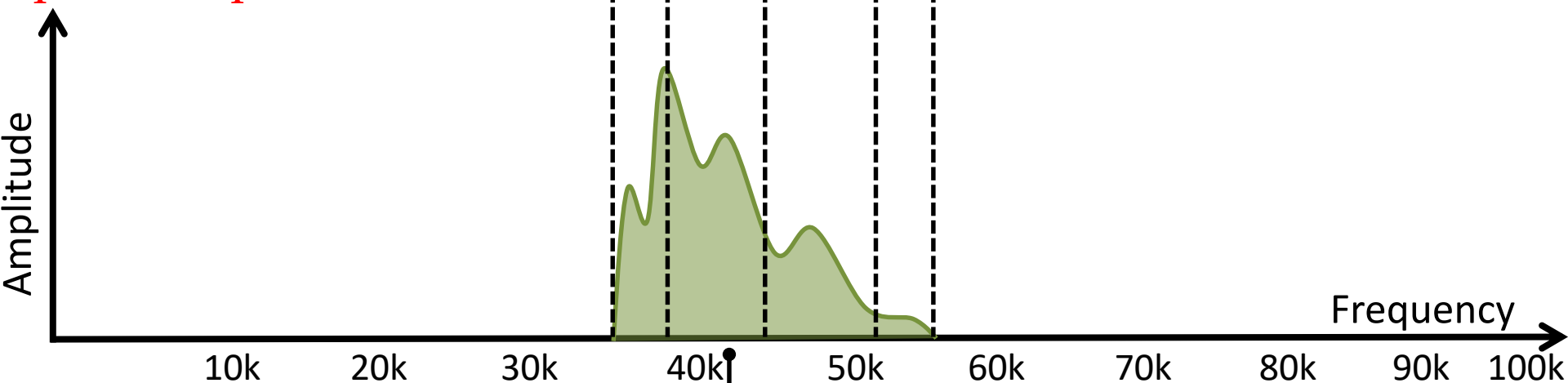


Speaker output

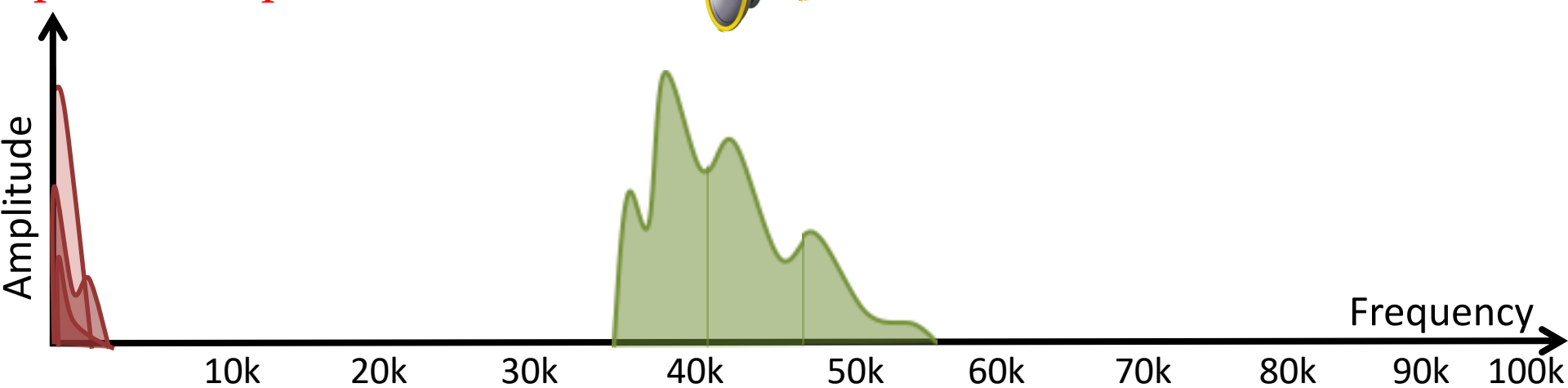


Speaker Nonlinearity \rightarrow Audible Leakage

Speaker input

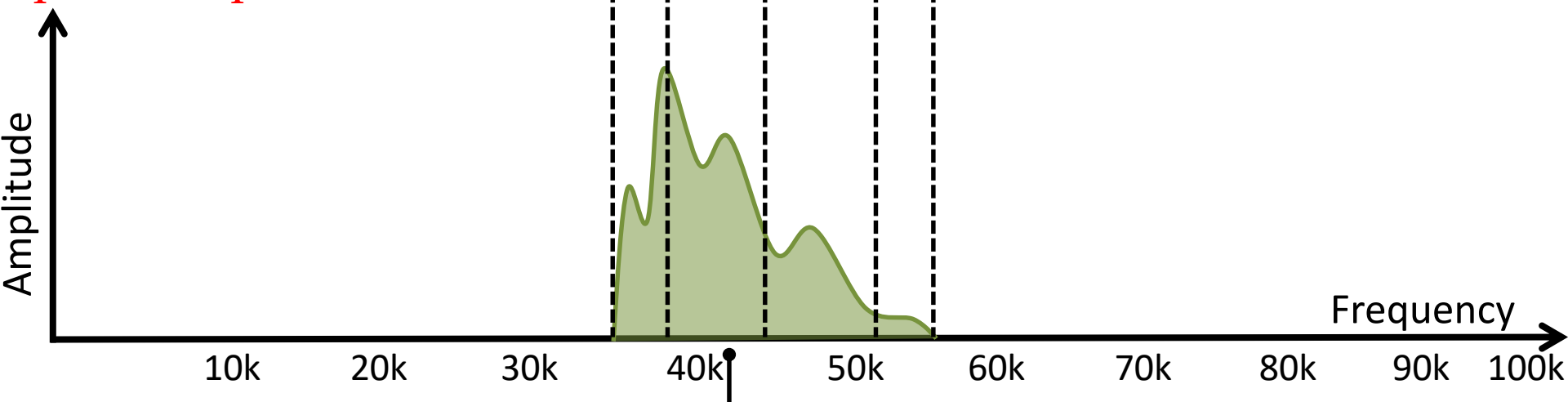


Speaker output

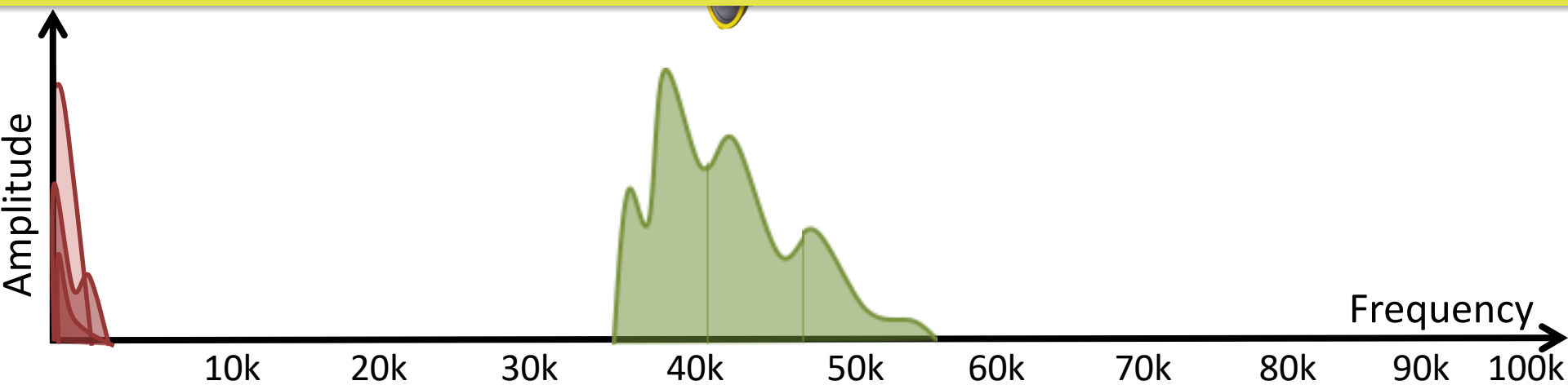


Speaker Nonlinearity \rightarrow Audible Leakage

Speaker input



Chopping compresses the leakage band



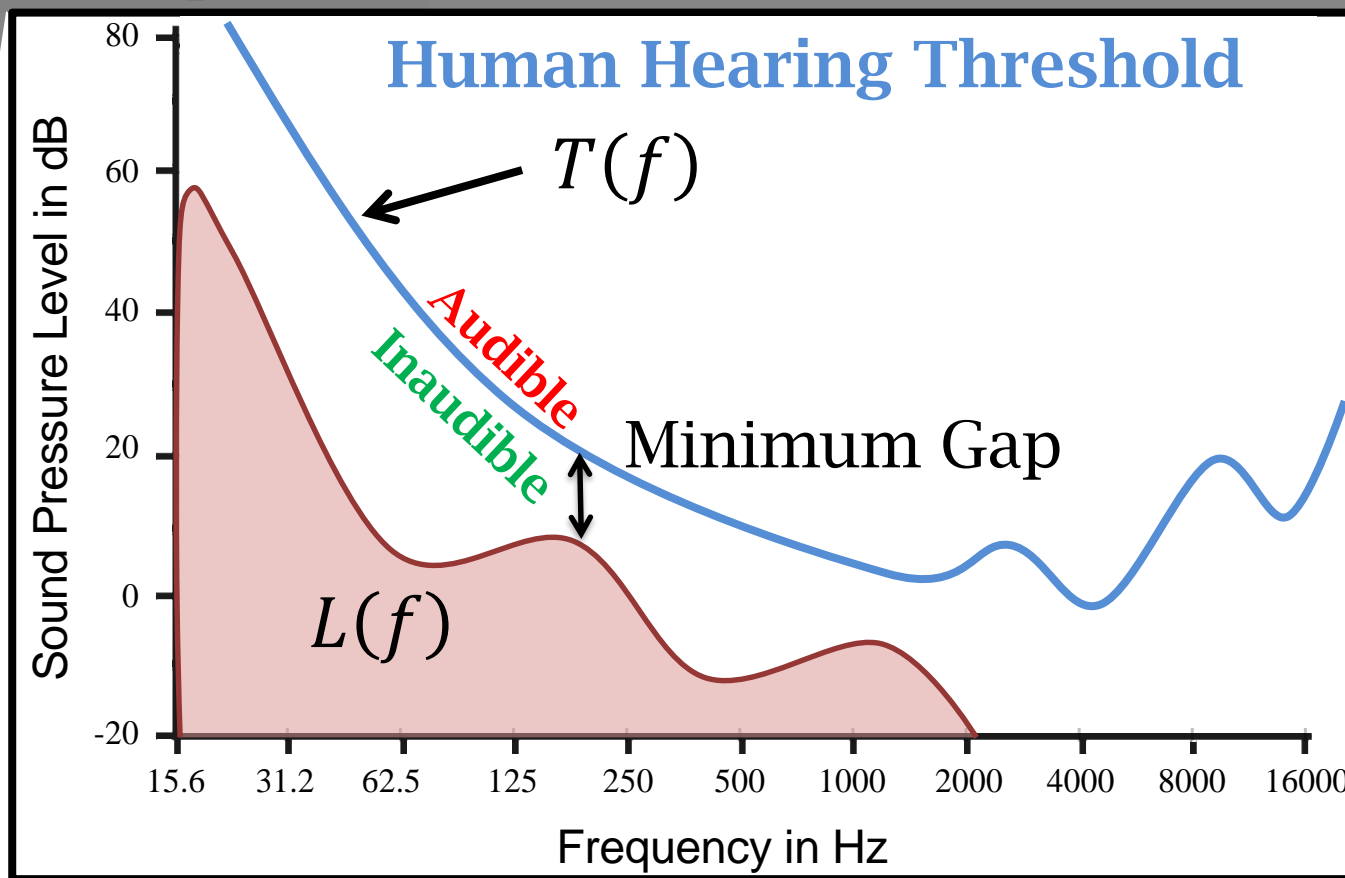
Leakage from Speaker

Speaker input

Amplitude

Speaker

Amplitude

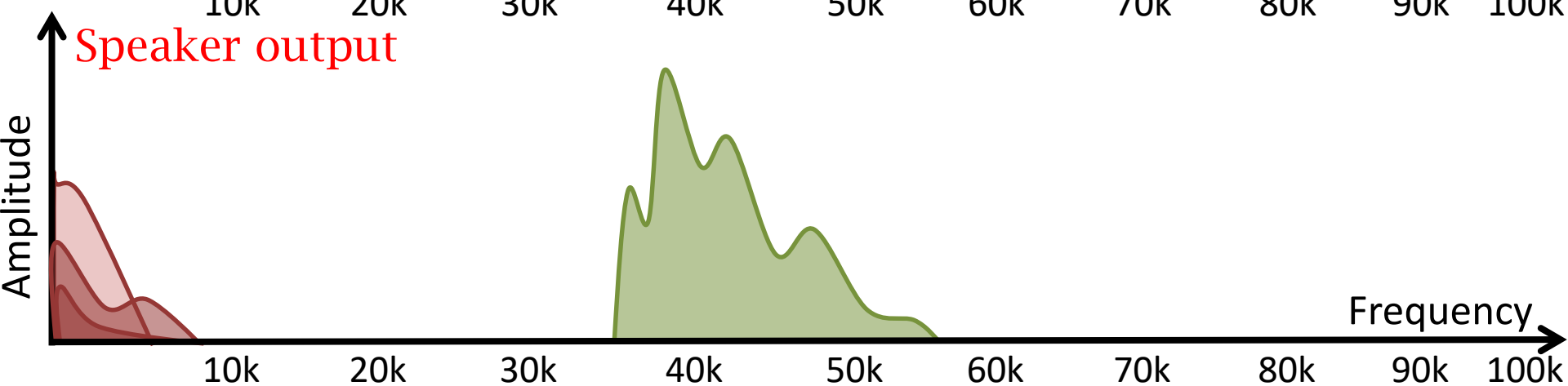
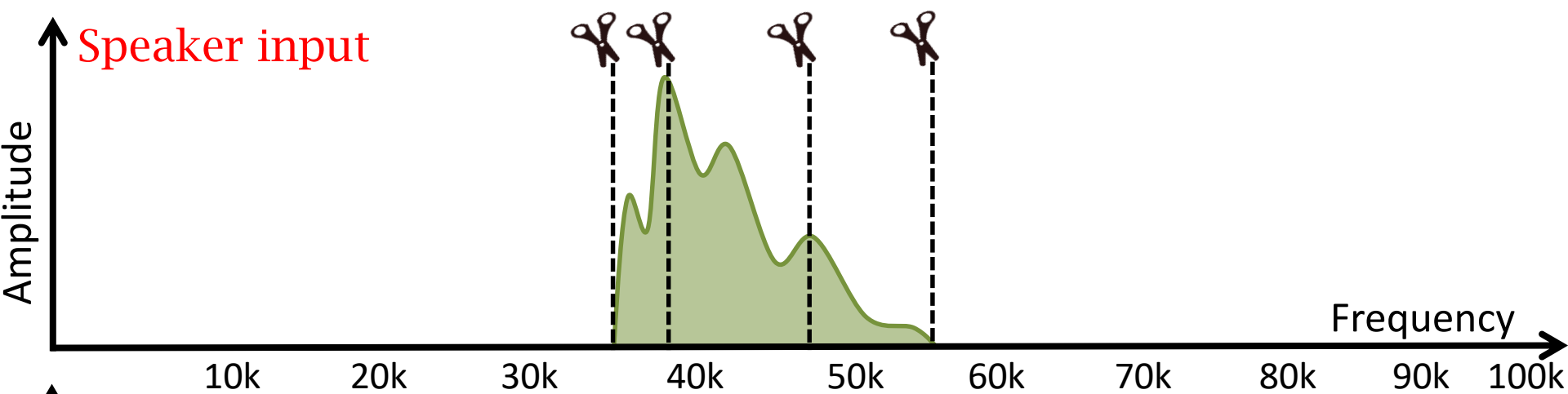


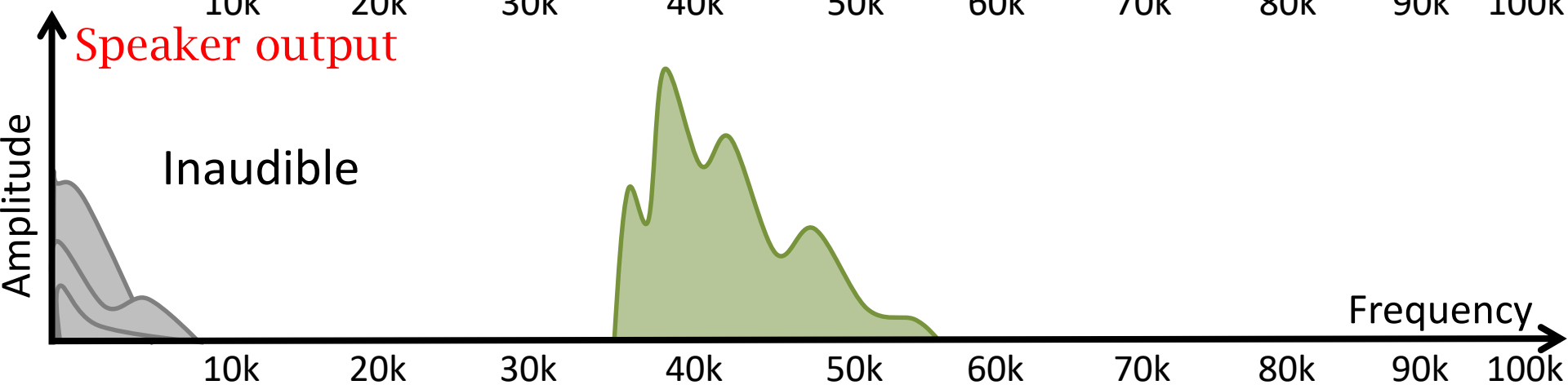
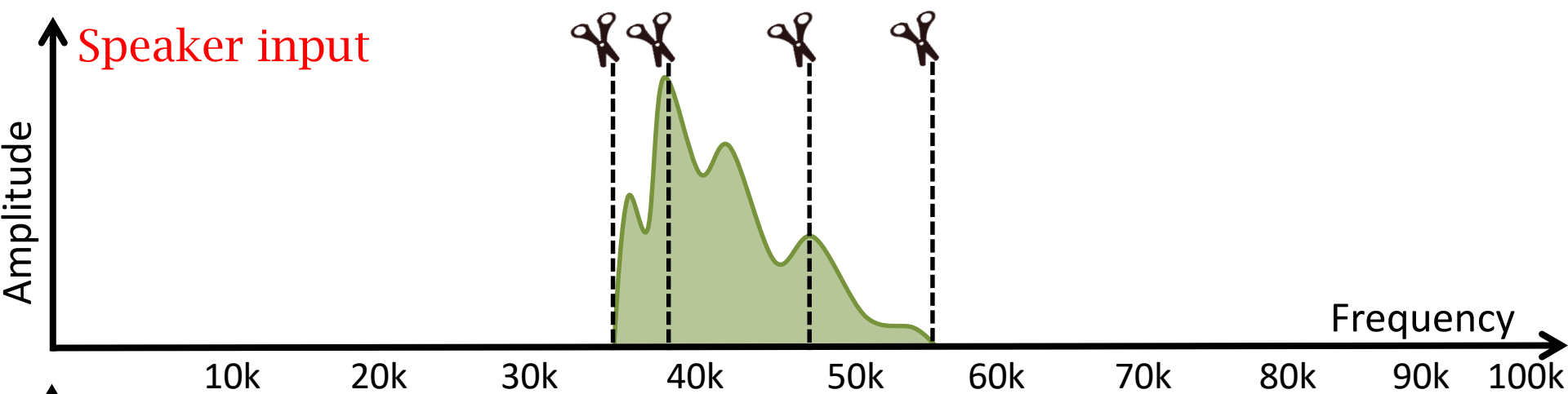
$$\text{Maximize } \min_f [T(f) - L(f)]$$

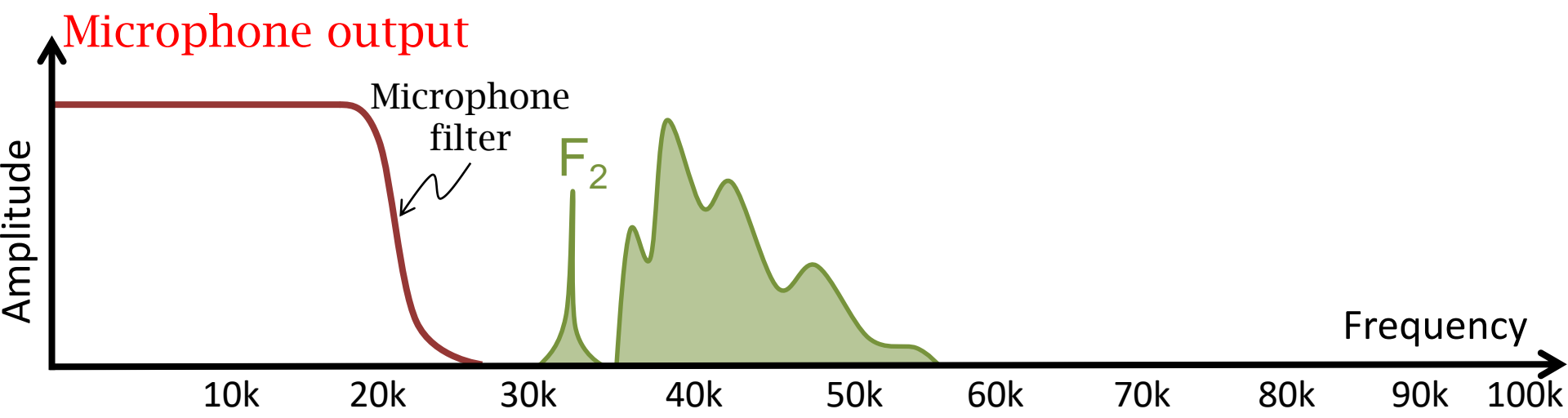
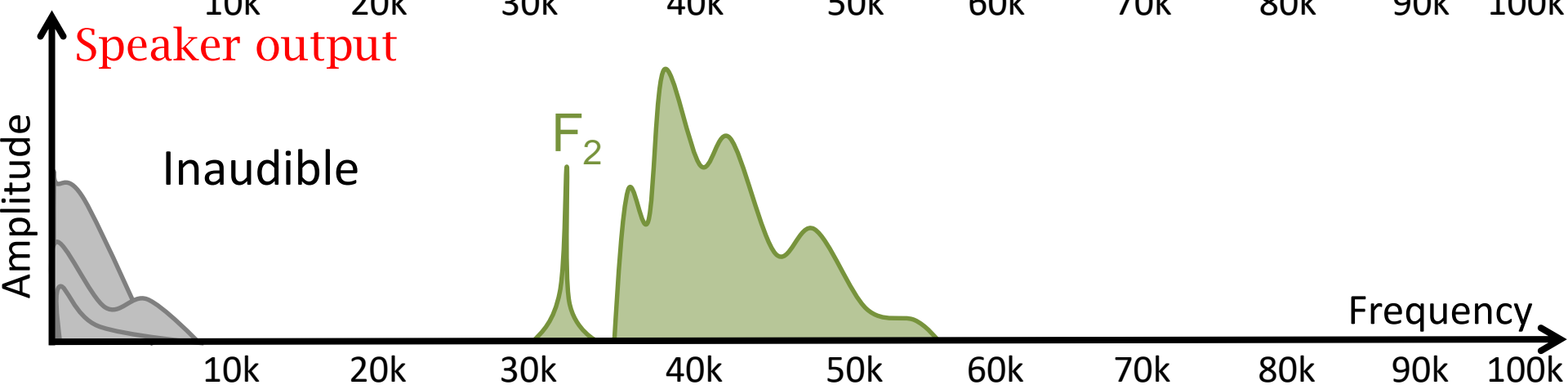
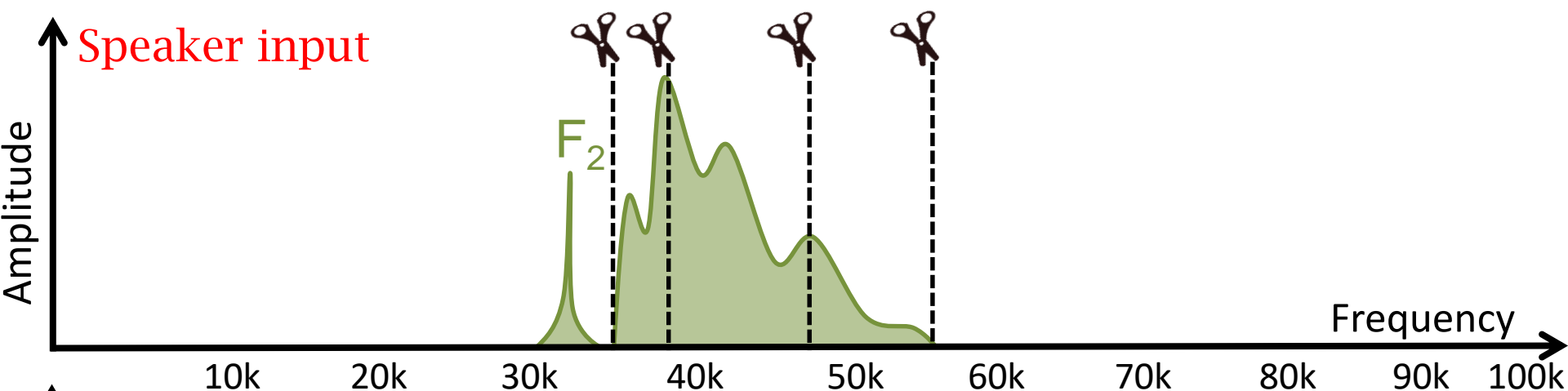
$$\text{subject to } f_0 \leq f_1 \leq f_2 \leq \dots \leq f_N$$

Frequency
90k 100k

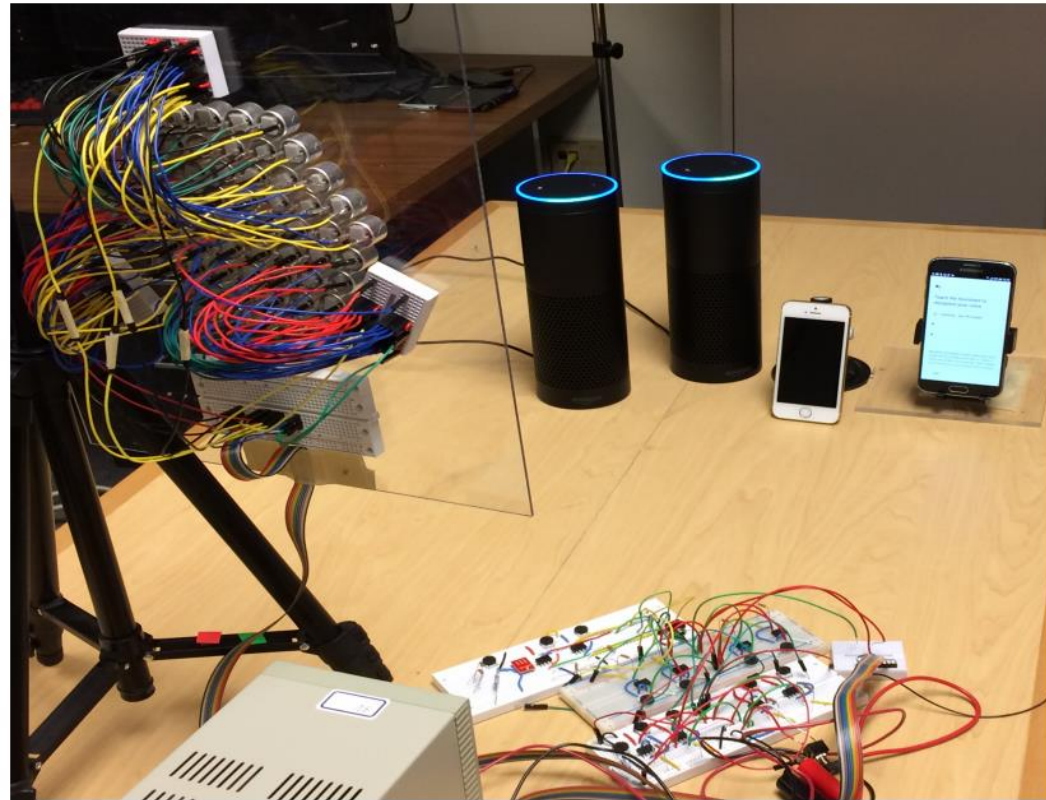
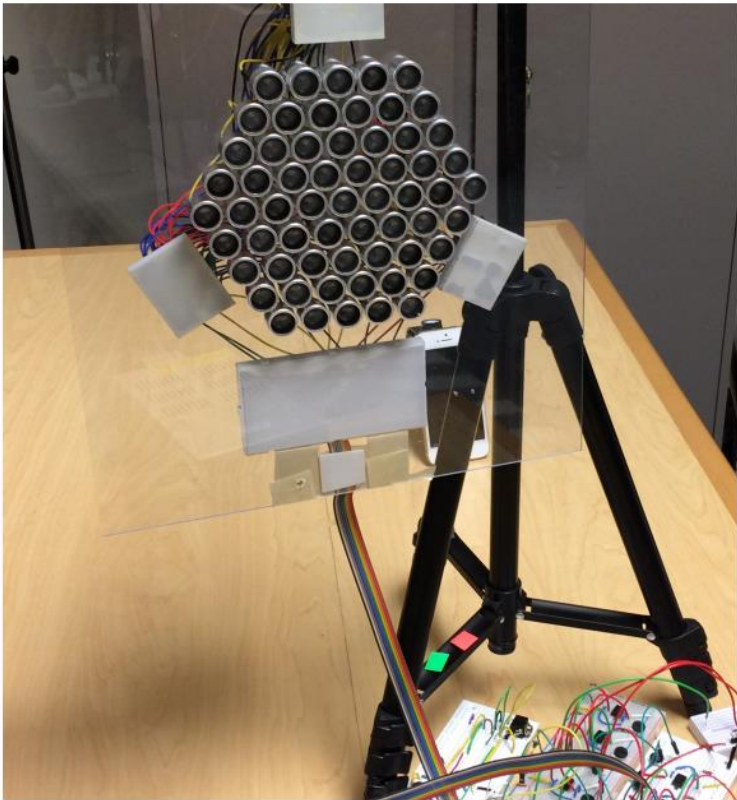
Frequency
10k 20k 30k 40k 50k 60k 70k 80k 90k 100k







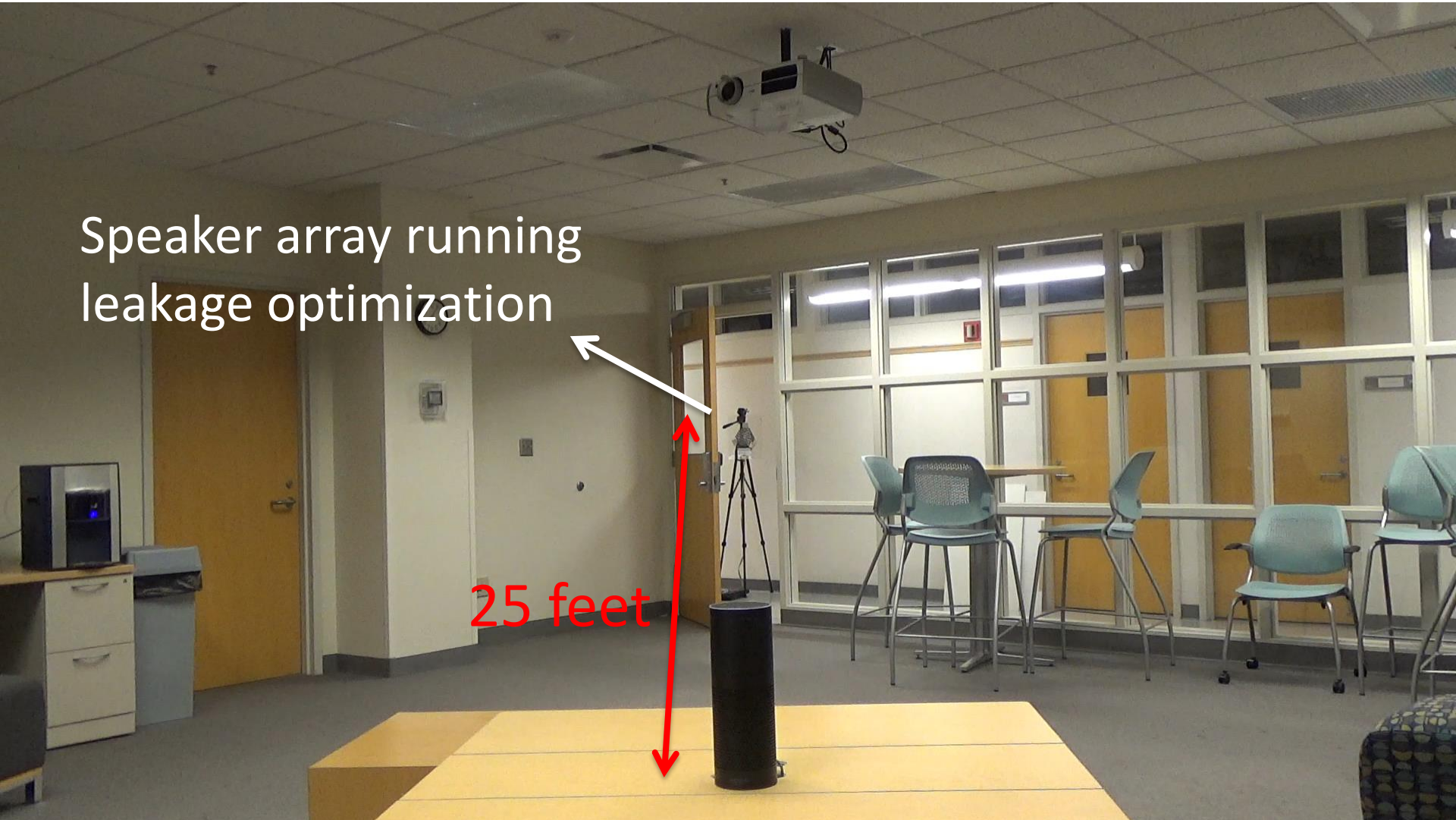
Evaluation



Inaudible voice commands: Long range

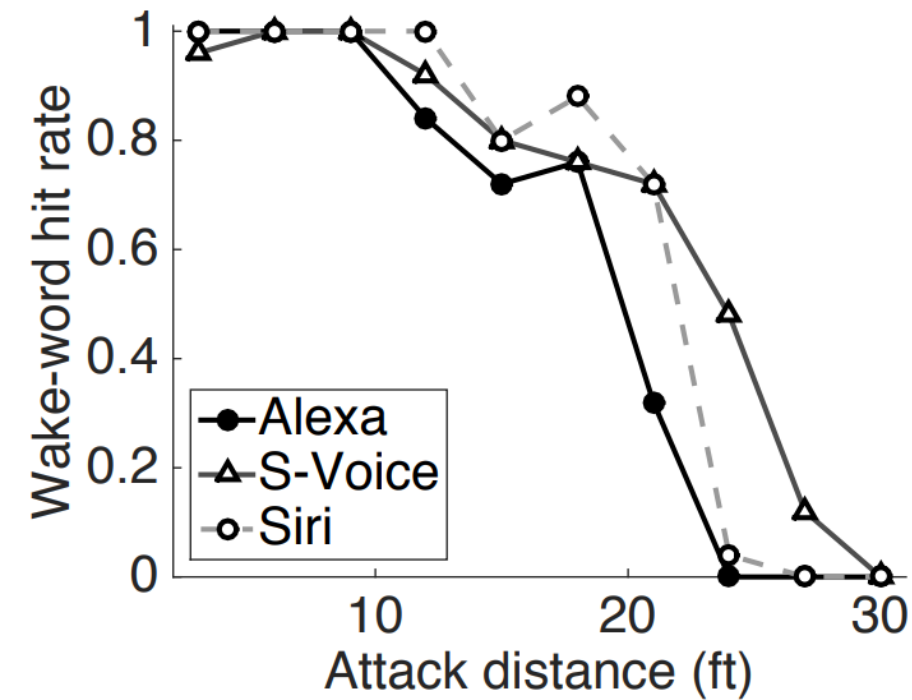
Speaker array running
leakage optimization

25 feet



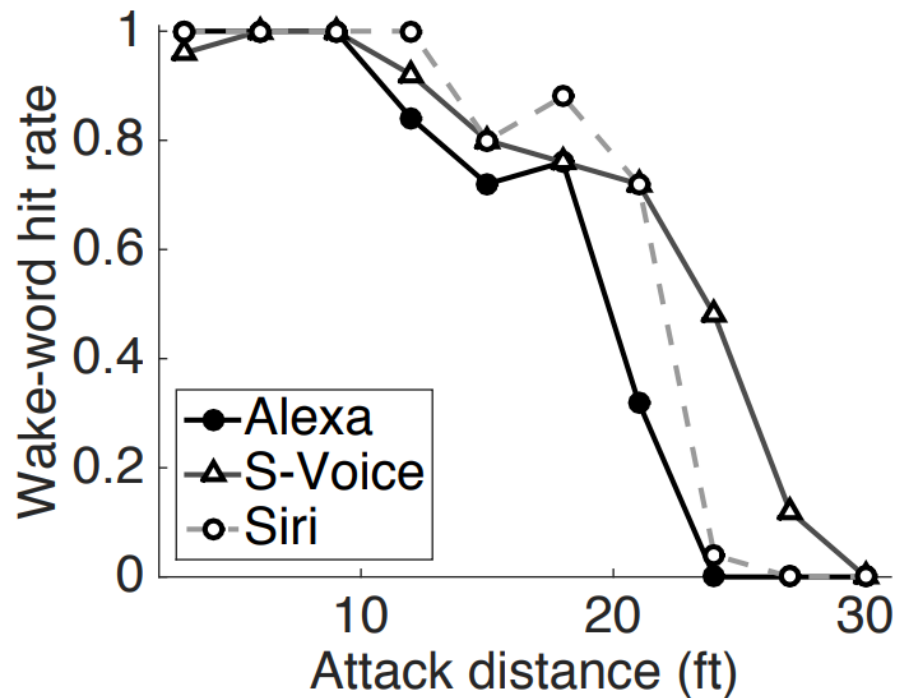
Evaluation

Wake-word hit rate

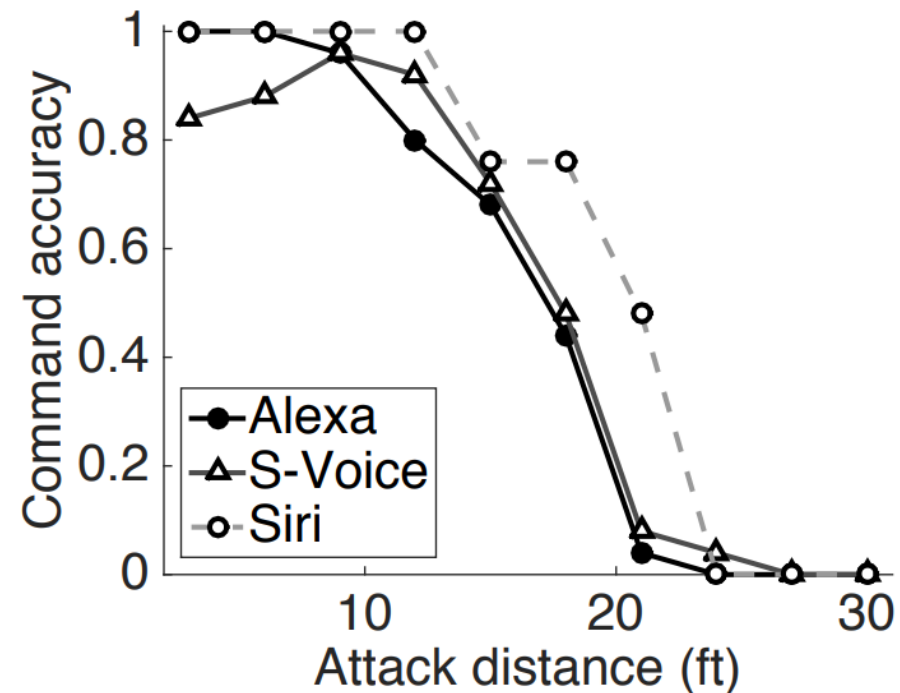


Evaluation

Wake-word hit rate

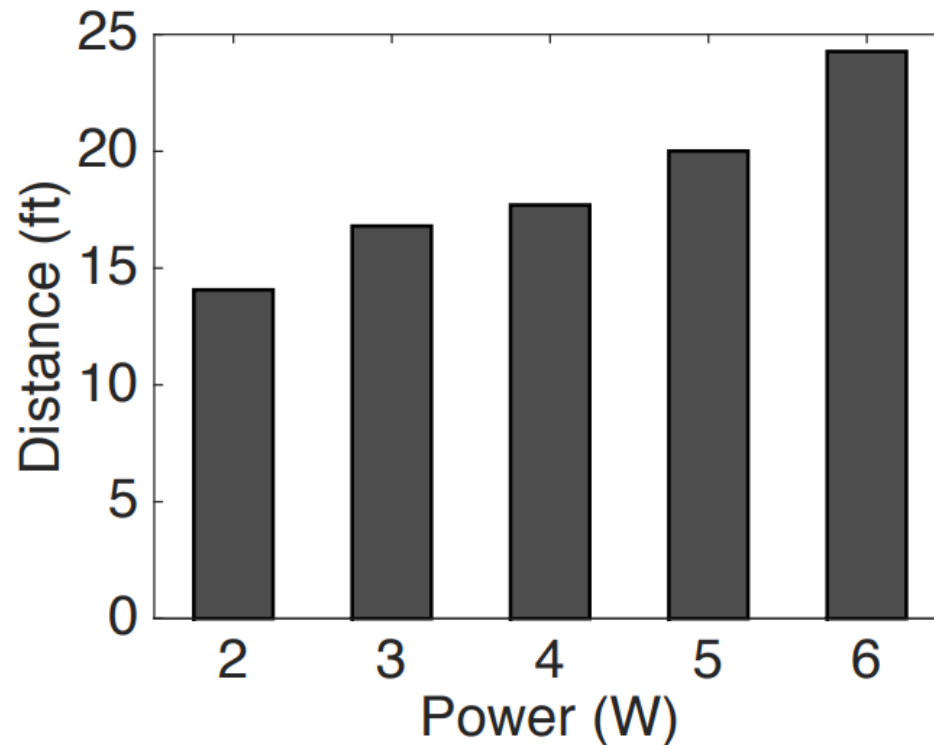


Command detection accuracy



Evaluation

Maximum activation distance for different input power



Talk Outline

0. [BackDoor], [DolphinAttack], [Princeton Video]

Today's Talk:

1. How to launch long-range (realistic) attacks?

2. How to defend against these attacks?

Talk Outline

0. [BackDoor], [DolphinAttack], [Princeton Video]

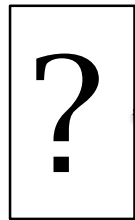
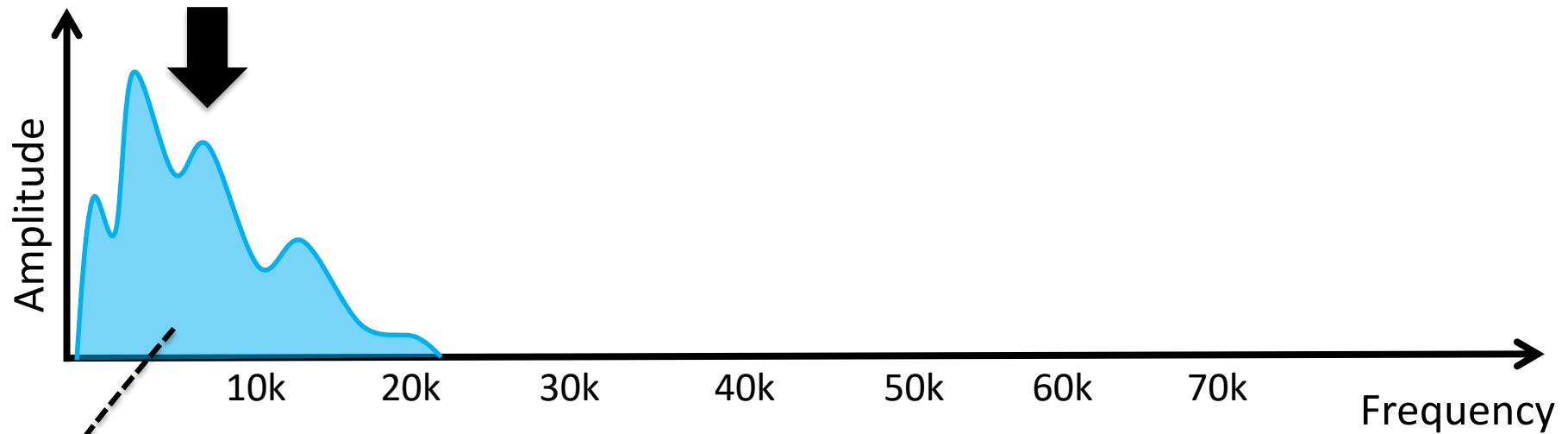
Today's Talk:

1. How to launch long-range (realistic) attacks?

2. How to defend against these attacks?

Core Question:

Is this a “non-linear signal” or normally recorded signal



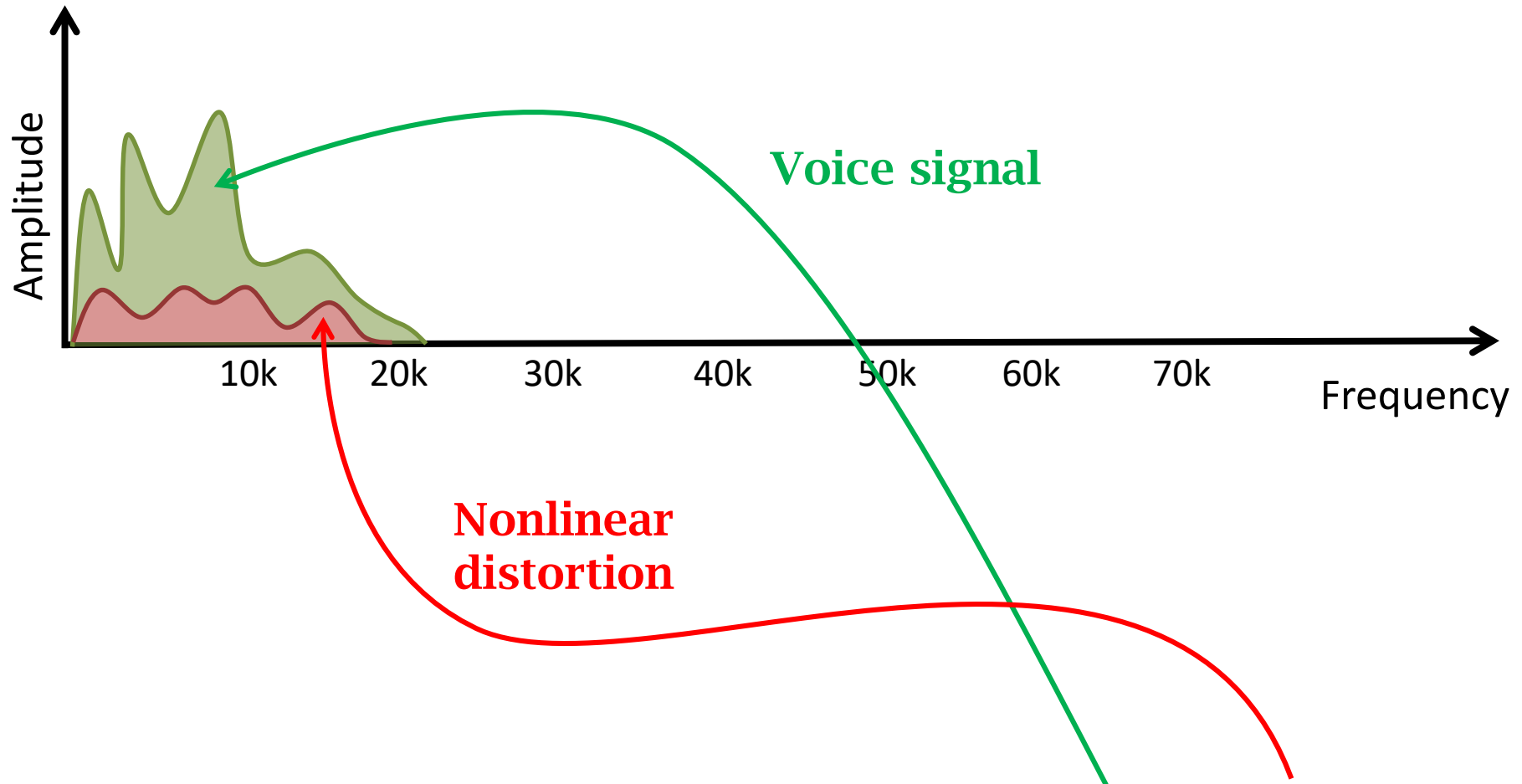
Voice signal $v(t)$

Inaudible
Voice Attack

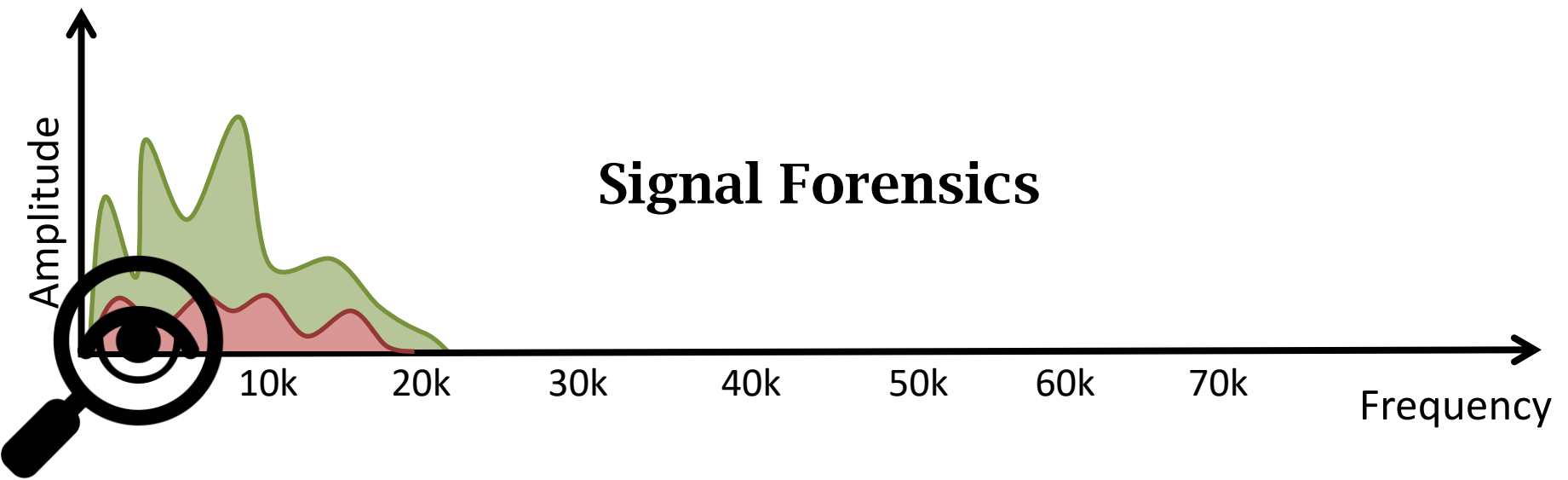
$$[v(t) \cdot \sin(\omega_1 t) + c \cdot \sin(\omega_1 t)]^2 \\ = v(t) + c'v^2(t) + \dots$$

Core Question:

Is this a “non-linear signal” or normally recorded signal



$$[v(t) \cdot \sin(\omega_1 t) + c \cdot \sin(\omega_1 t)]^2 = v(t) + c'v^2(t) + \dots$$

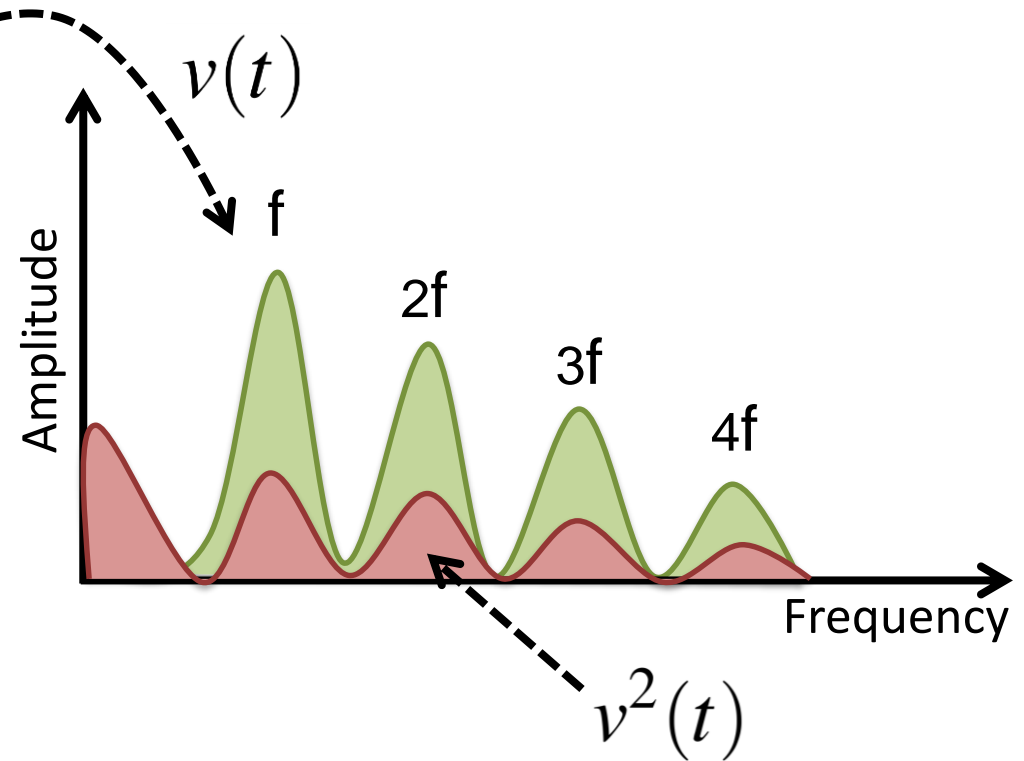
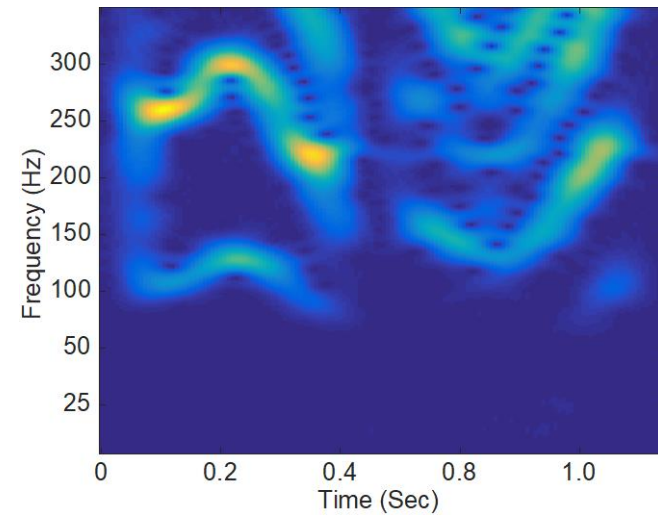


Difficult to decouple “voice signal” and “non-linear signal”

Human voice signals present opportunities ...

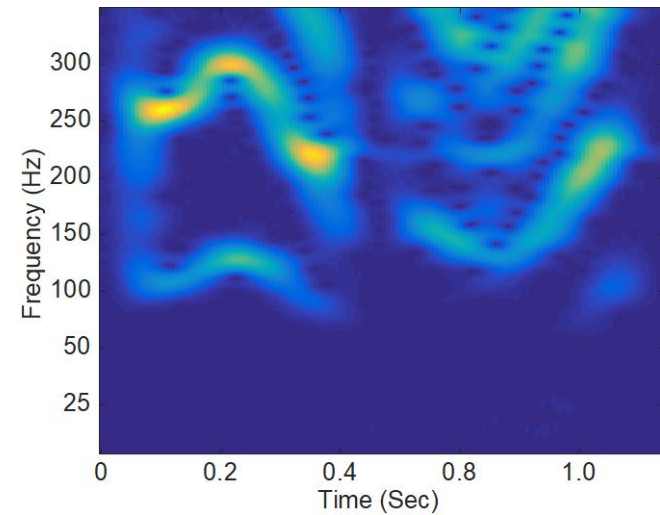
Opportunity #1: Voice > 50 Hz

Human Voice

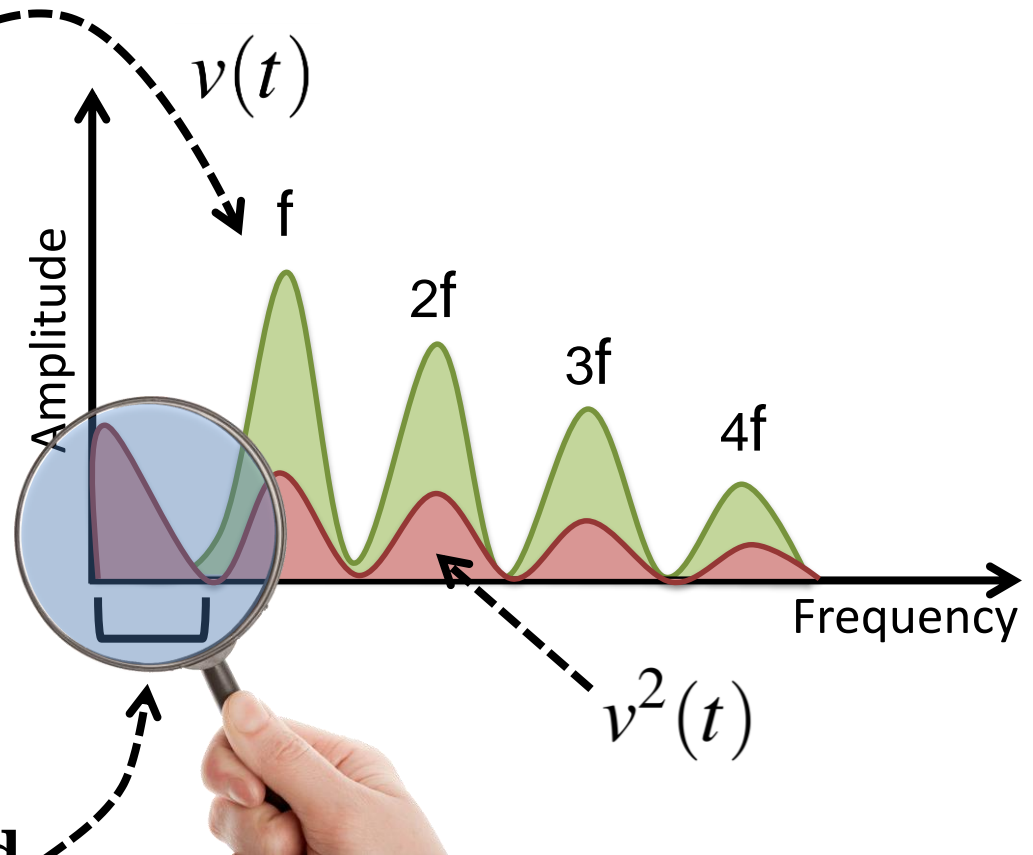


Opportunity #1: Voice > 50 Hz

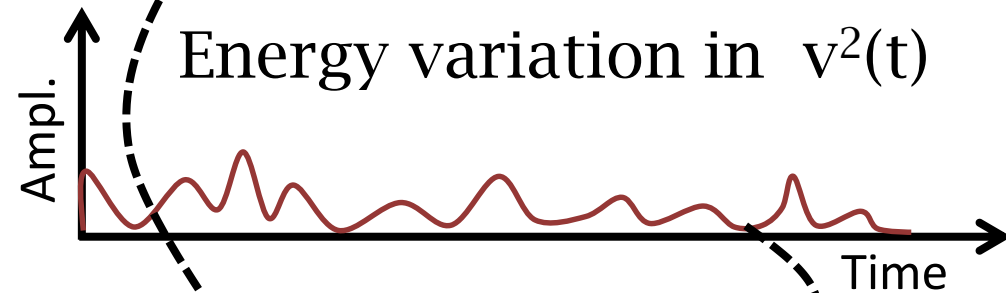
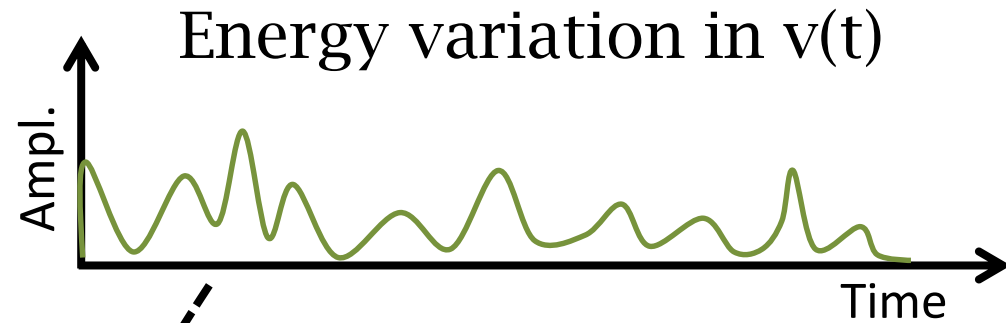
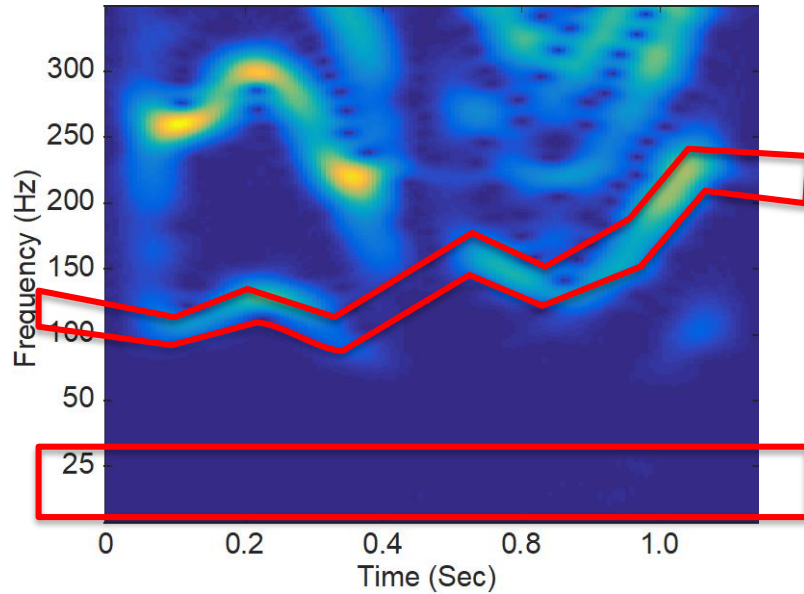
Human Voice



Energy at
sub-50Hz band

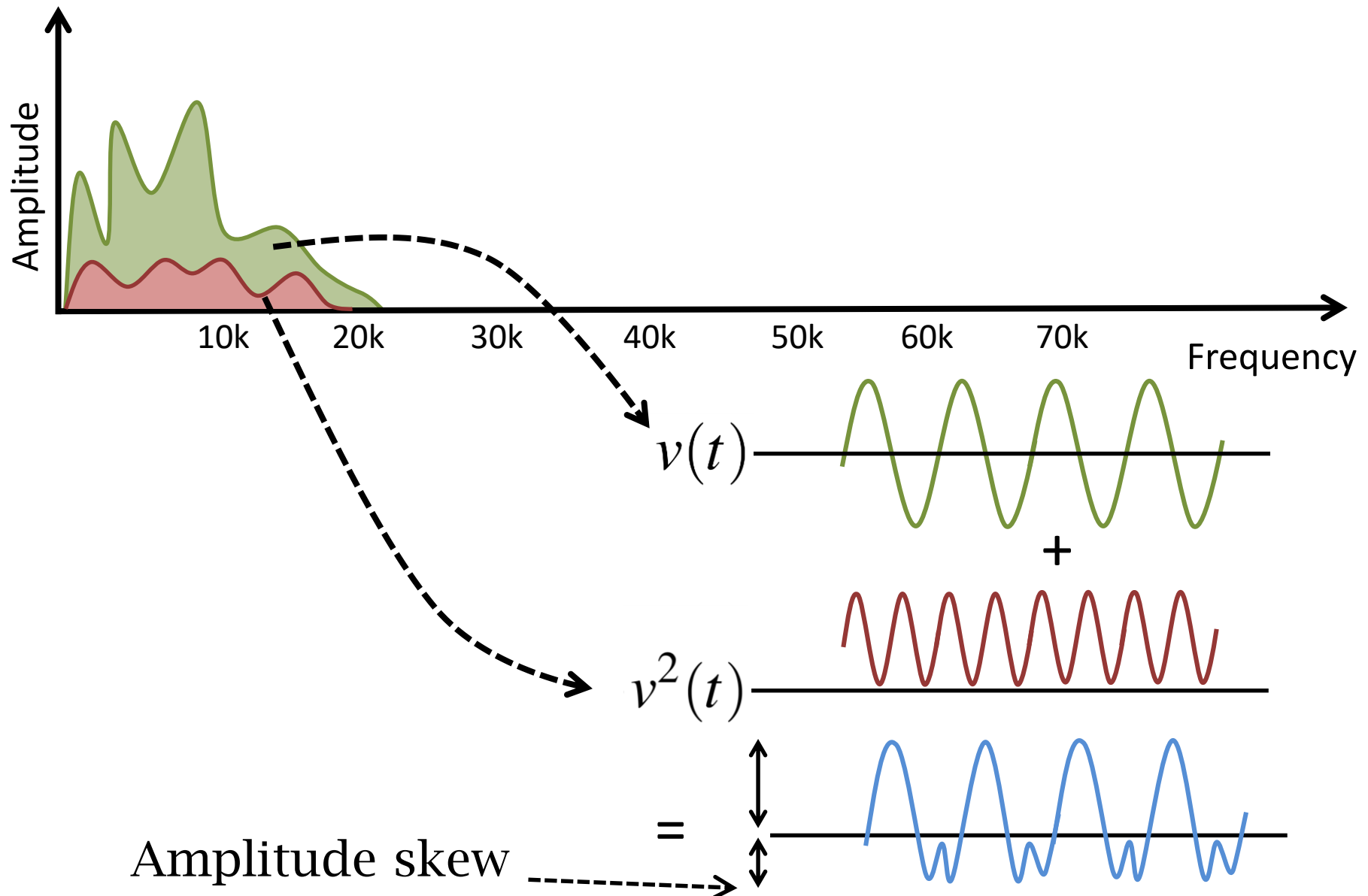


Opportunity #2: Correlation

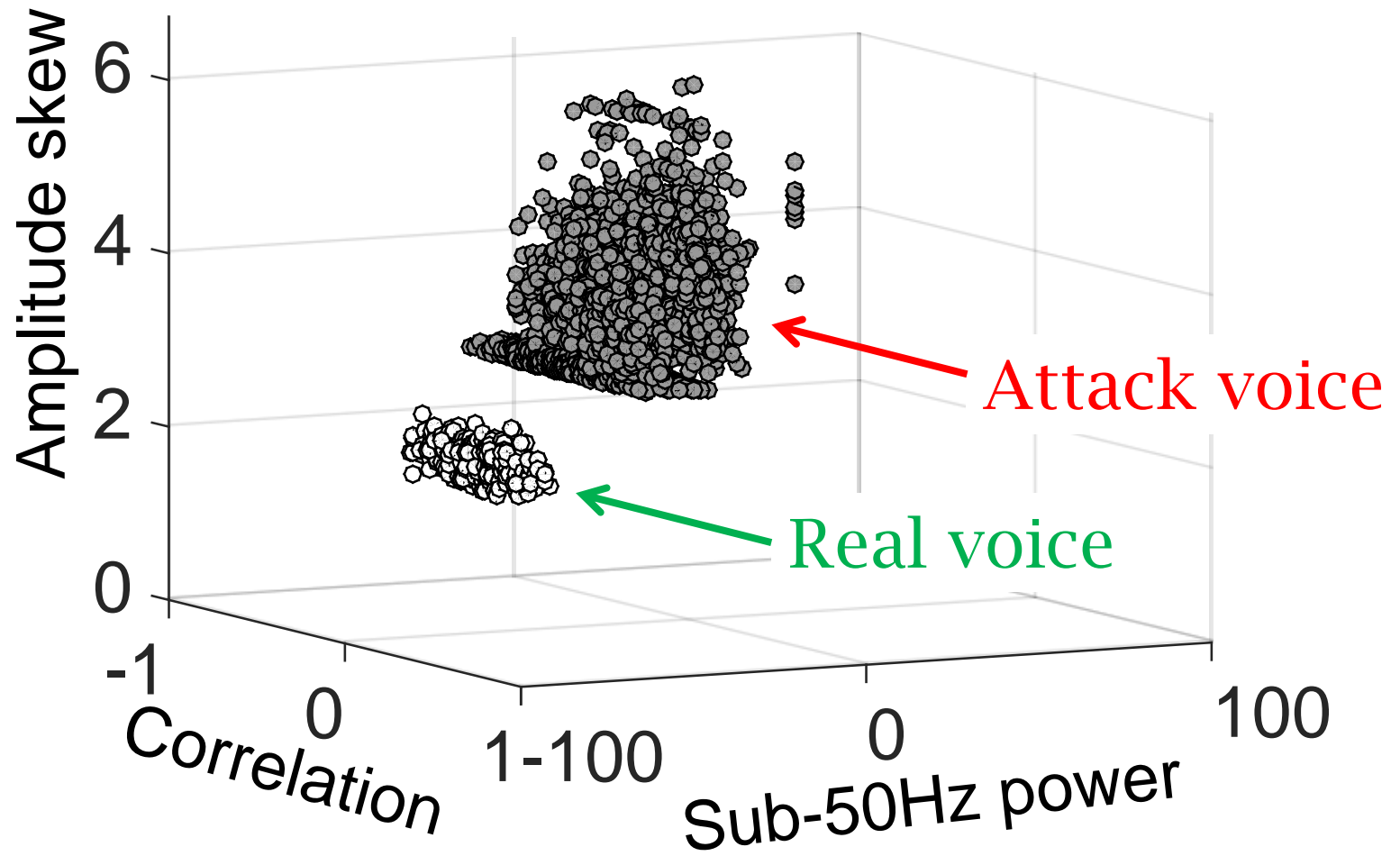


Correlation

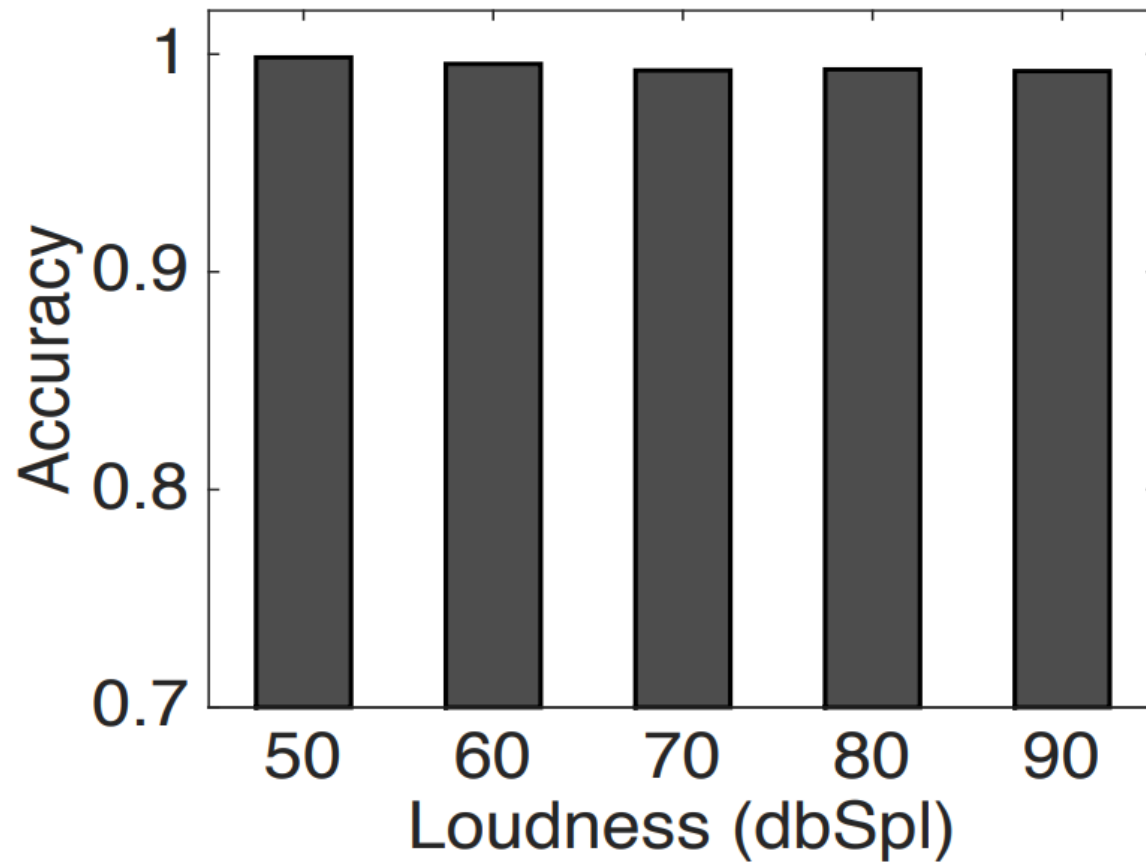
Opportunity #3: Amplitude Skewness



5000 Test Cases



Overall Detection Accuracy



To summarize...

Inaudible Acoustics (> 25 kHz): “Alexa, open the garage door!”

