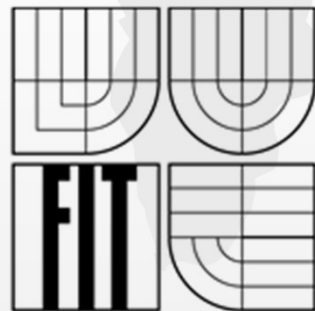


ZRE - Kódování řeči II.

CELP

Vladimír Malenovský, ÚPGM FIT VUT Brno

malenov@fit.vutbr.cz



Úvod

- CELP vznikl jako pokus o zlepšení kvality LPC kodéru
- CELP je jedna z nejužších myšlenek v oblasti kódování řeči
- Akronym CELP vymysleli Manfred Schroeder a Bishnu Atal v roce 1985 v článku

Schroeder, M.R. and B. S. Atal (1985). "Code-Excited Linear Prediction (CELP): High-Quality Speech at Very Low Bit Rates," IEEE ICASSP, pp. 2511–2514.



Úvod

- nedostatky kodeku LPC
 - příliš zjednodušené kódování buzení - pouze 2.4 kbps
 - striktní rozlišení znělé a neznělé řeči
 - fáze vstupního signálu není zachována
- hlavní „nové“ myšlenky CELPu:
 - **dlouhodobý prediktor (LTP)** – znělou řeč lze krásně předpovědět i ze vzdálené minulosti
 - **kódování obou složek buzení naráz, znělé i neznělé** – lidé přece neříkají jen „a“ a „s“
 - **koncept analýzy syntézou (analysis-by-synthesis approach)** – učení se z vlastní chyby
 - **perceptuální váhování** – co ucho neslyší, to kodek nekóduje

- nové kodeky:

FS 1016 (CELP)

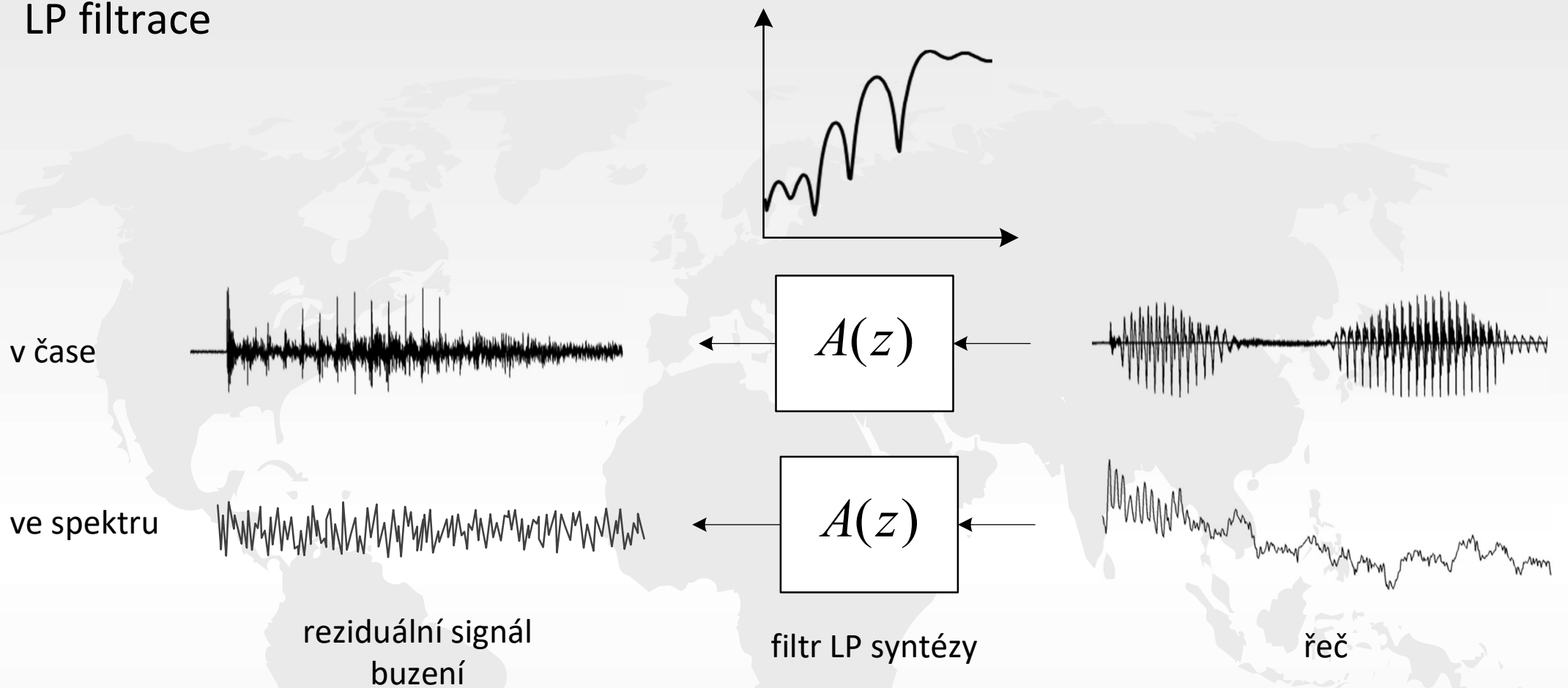
GSM-FR (RPE LTP)

GSM-HR (VSELP)



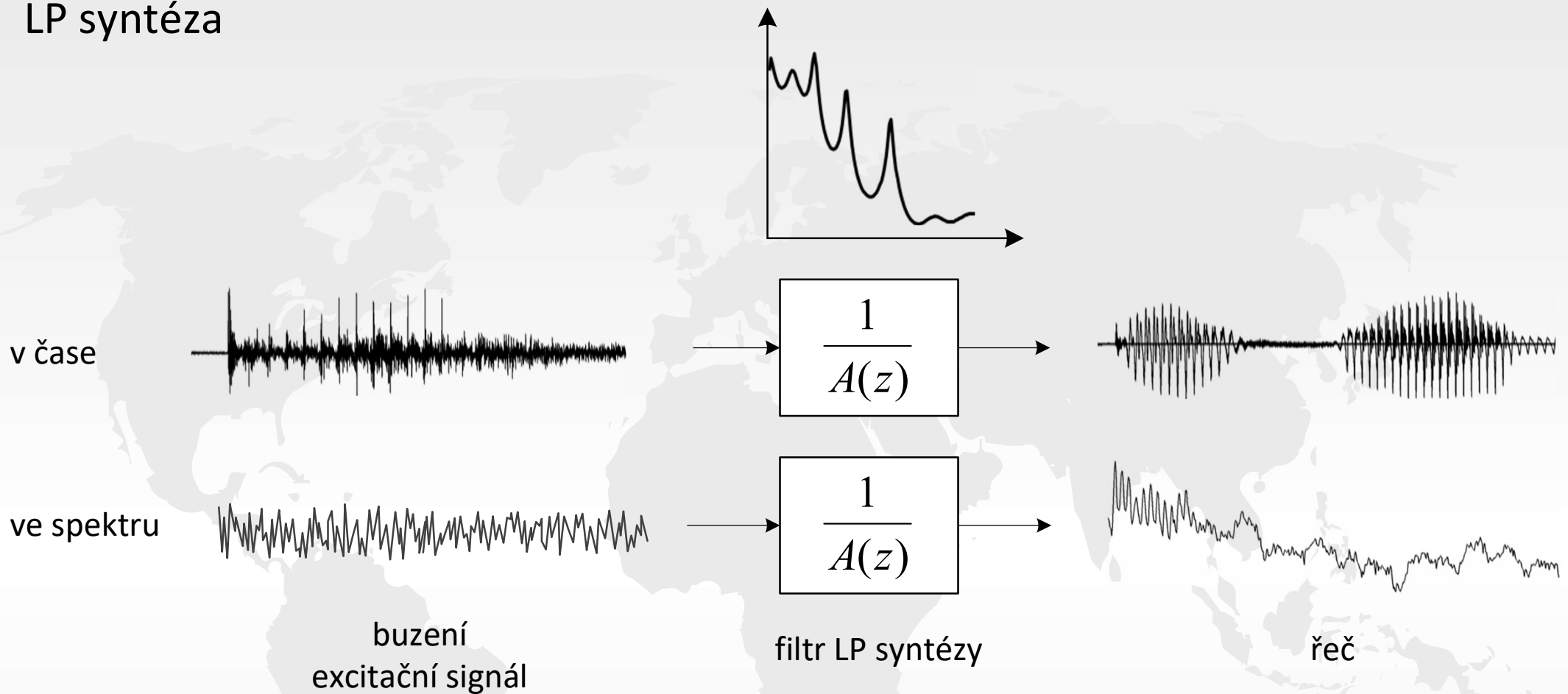
LP model tvorby řeči

LP filtrace



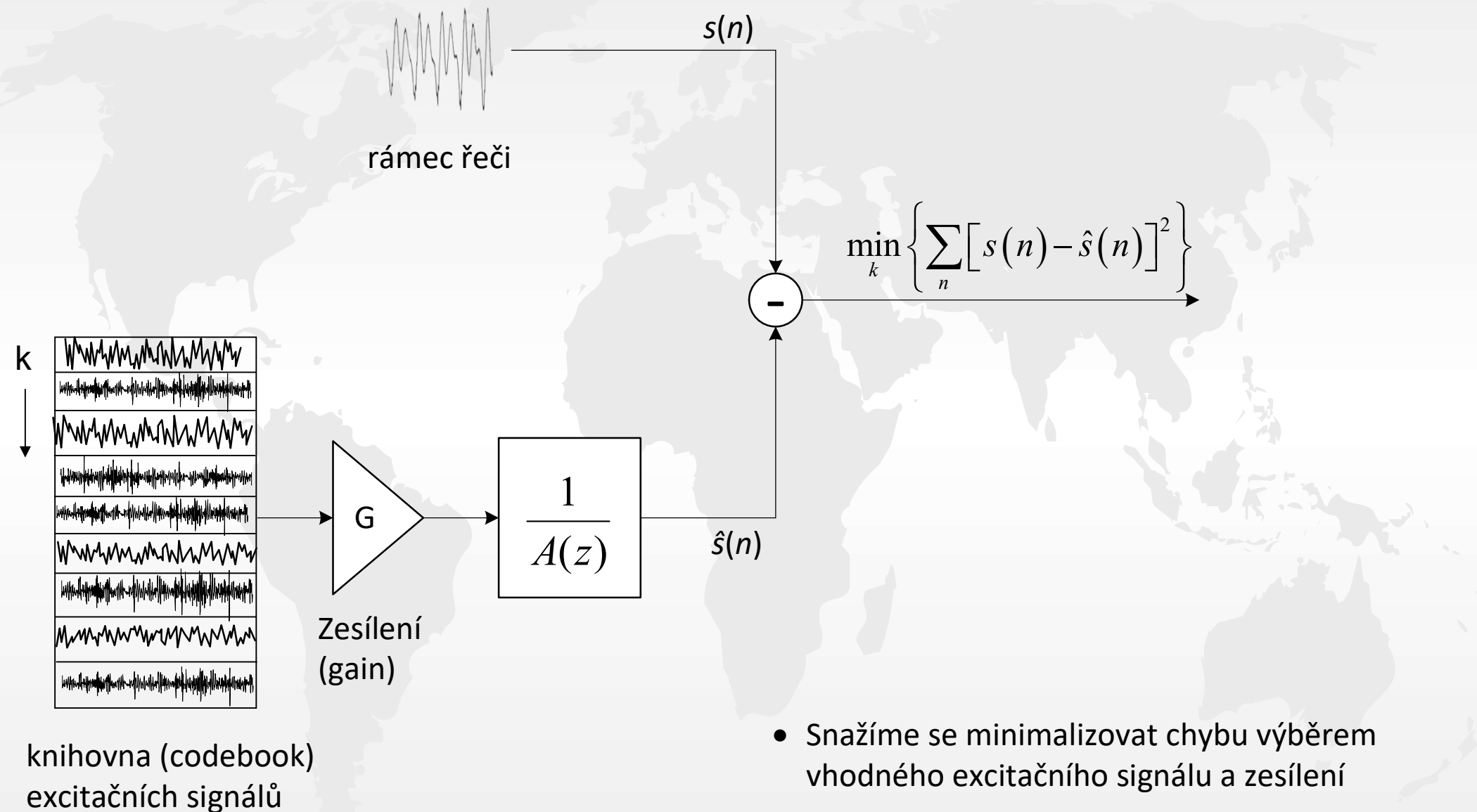
- Filtrací vstupního signálu filterm $A(z)$ odstraňuje z řečového signálu jeho vlastní „obálku“ a tím zbavuje spektrum vlastních formantů
- Reziduální signál (buzení) je to co po této operaci zbyde a to se dále kóduje

LP syntéza



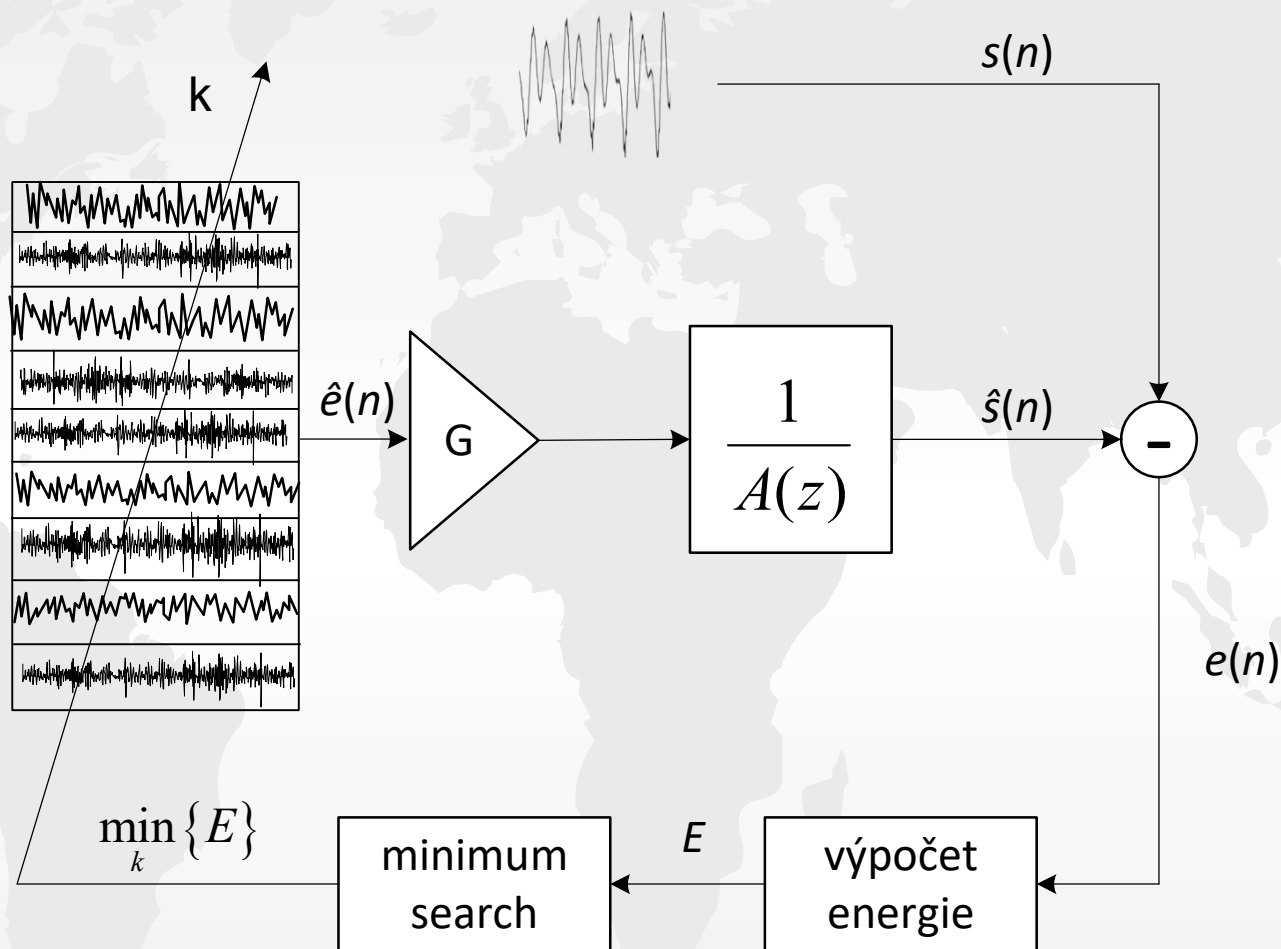
- téměř všechny řečové kodeky pro veřejné komunikace v dnešní době používají LP model
- filtr LP syntézy je filterm typu IIR, takže má vlastní paměť, která odpovídá několika posledním vzorkům z minulosti řeči
- filtr LP syntézy má pouze póly, t.j. umí modelovat pouze spektrální „špičky“, nikoliv „zářezy“
- vzhledem k charakteru filtru může dojít k jeho nestabilitě a „explozi“ syntézy

CELP – princip kodéru

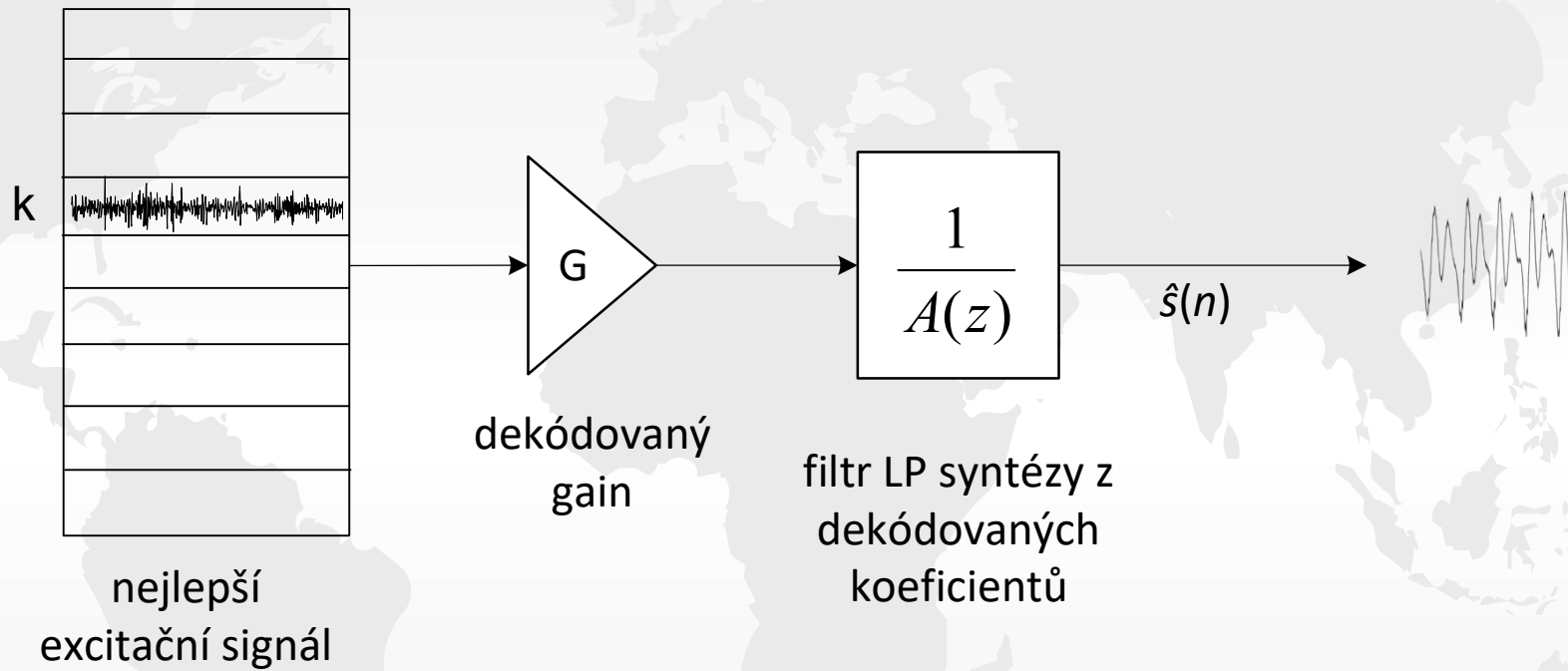


Koncept analýzy syntézou

překreslené schéma

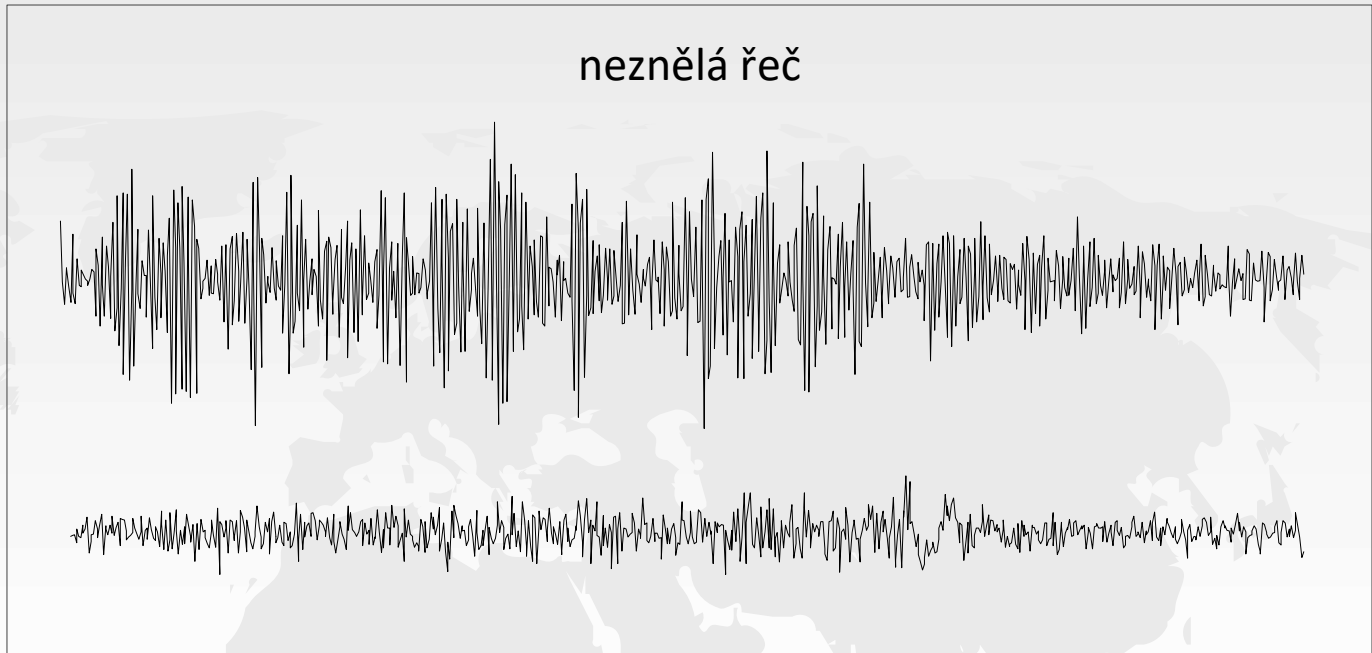


CELP – princip dekodéru

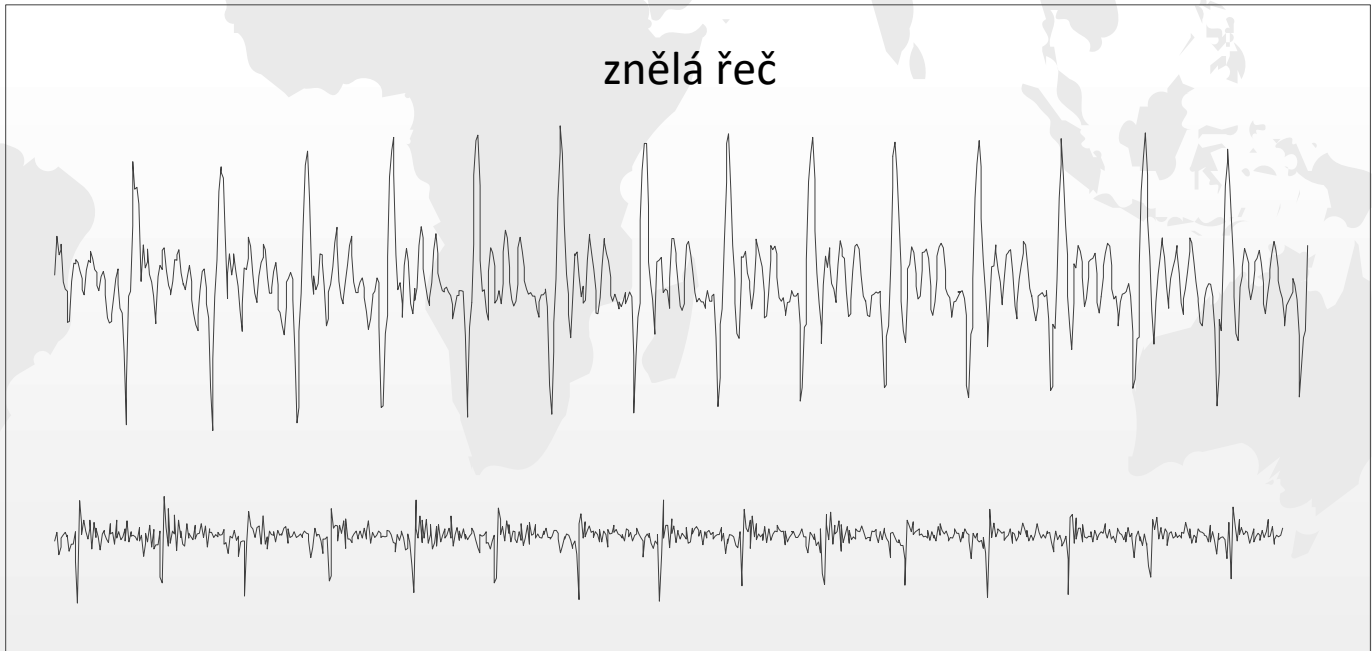


Reziduální signál

neznělá řeč

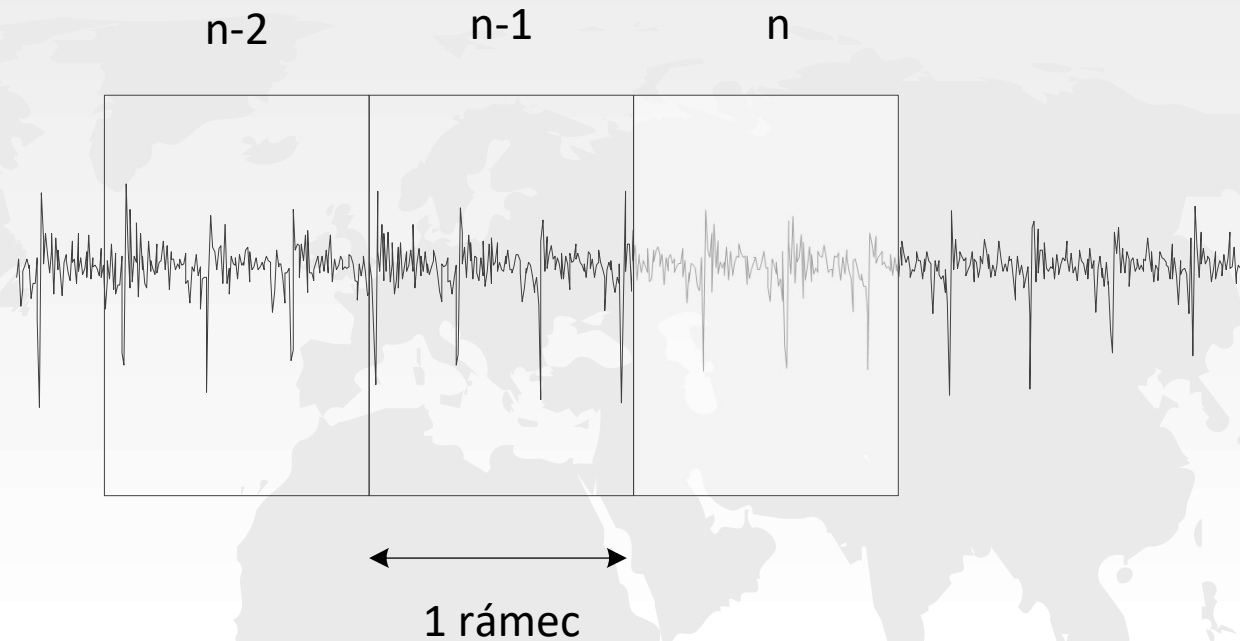


znělá řeč



periodicita v LP reziduu

Reziduální signál

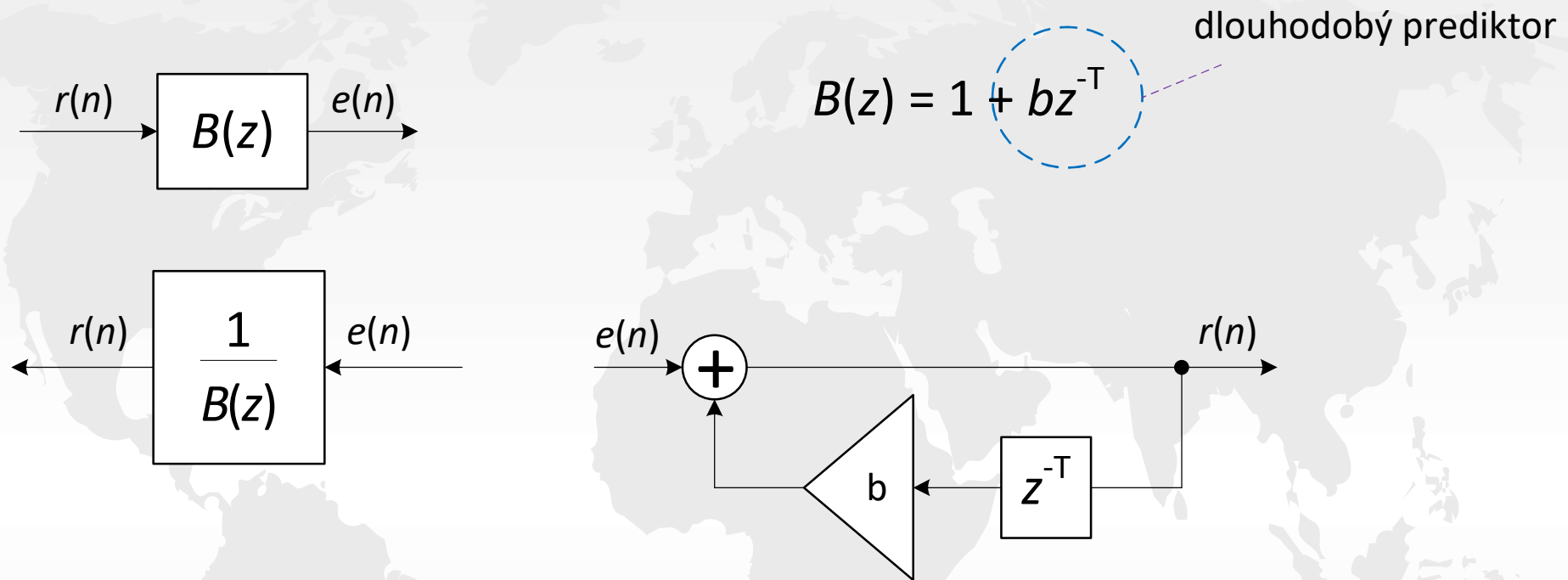


- reziduální signál v současném rámci lze predikovat z minulosti
- u znělého signálu se úseky „opakují“ v periodách odpovídajících délce základního tónu
- můžeme toho využít a kódovat pouze rozdíl mezi současným a predikovaným signálem
- ale pozor, k predikci musíme použít již zakódovaný (přenesený) rez. signál a ne originál, protože jinak by enkodér a dekodér nepracovali se stejnými signály



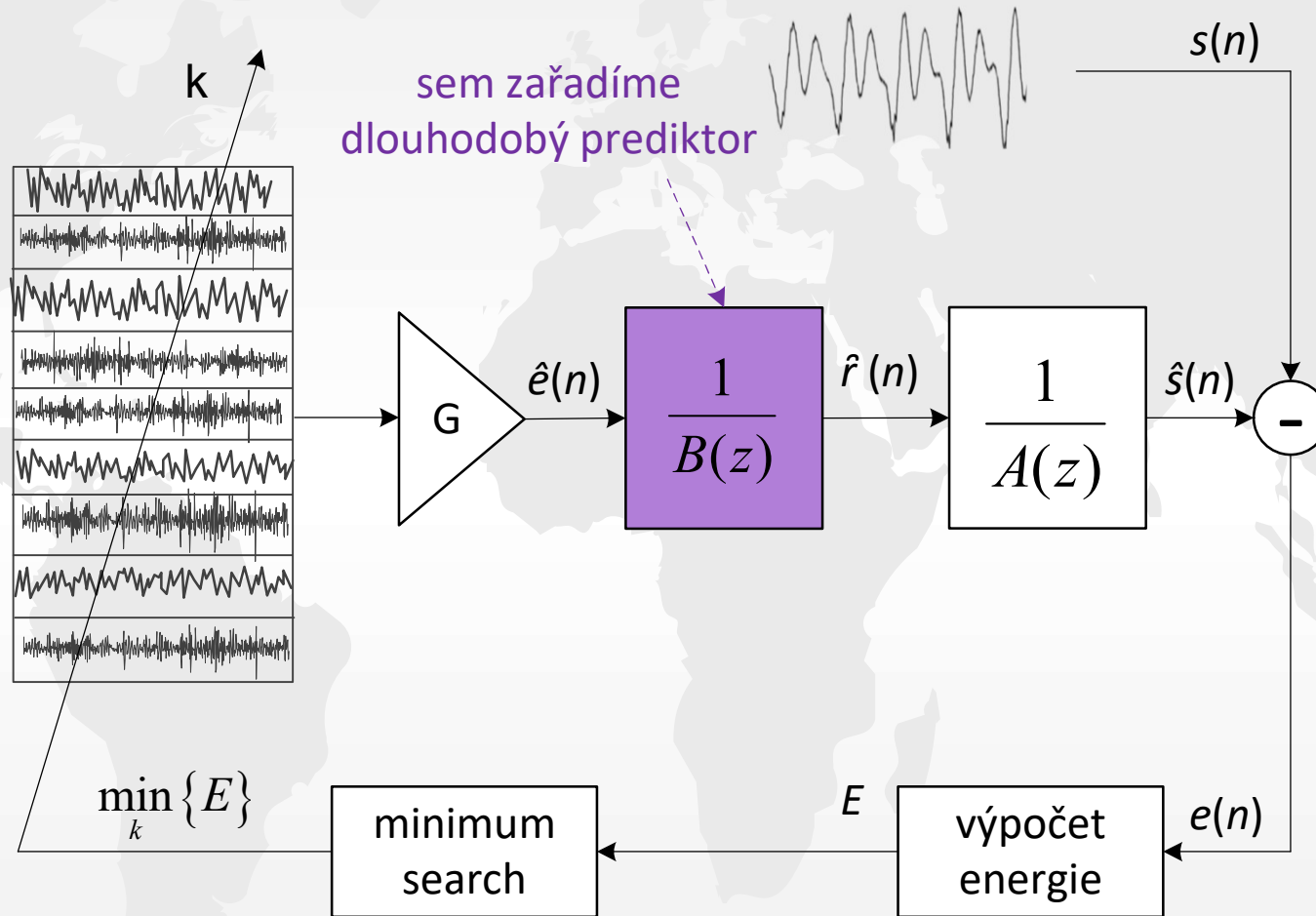
Dlouhodobý prediktor (LTP)

Dlouhodobý prediktor (Long-Term Predictor)

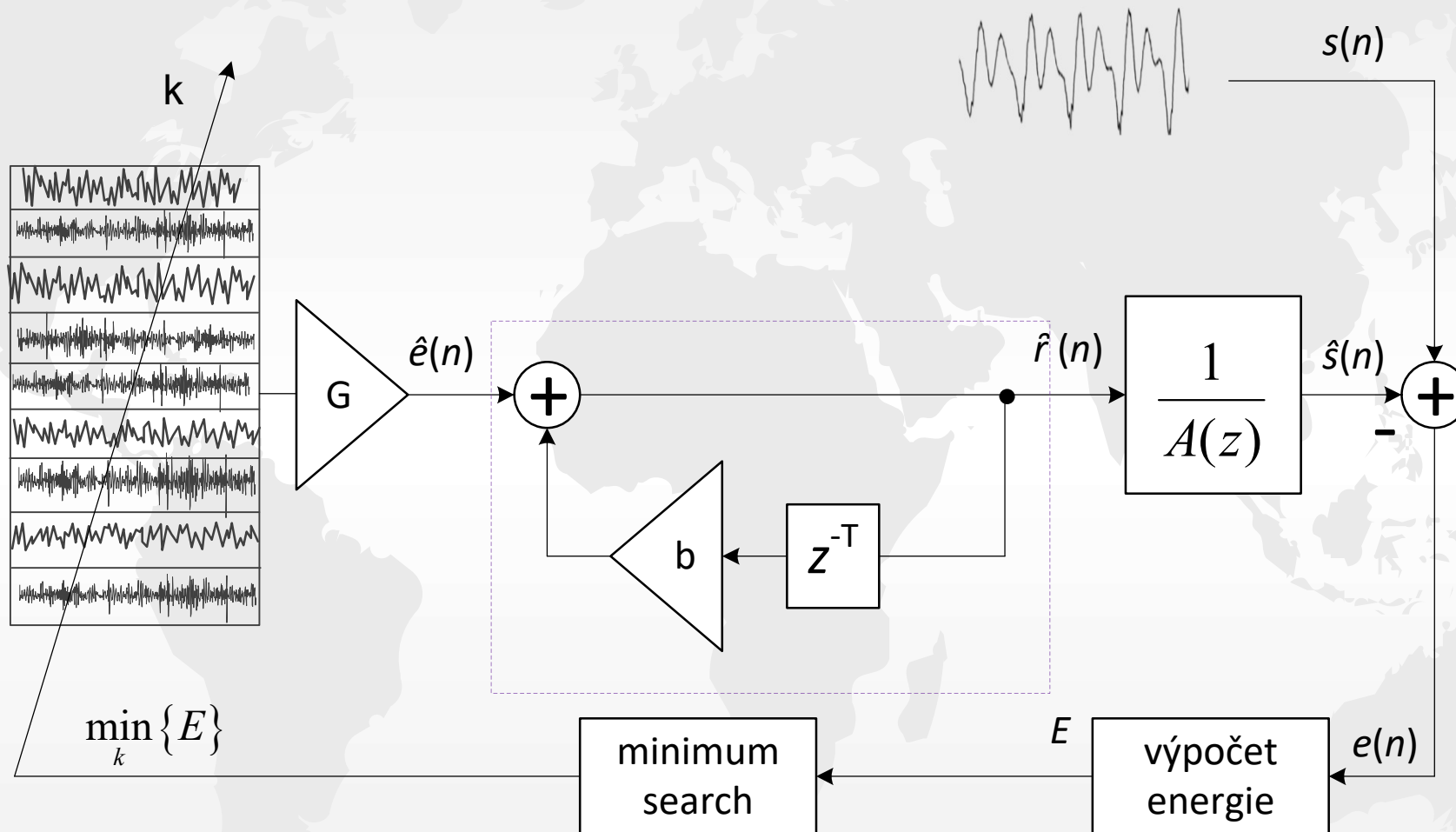


- dlouhodobý prediktor předpovídá současný vzorek signálu $r(n)$ ze minulého vzorku $r(n-T)$
- chybový signál: $e(n) = r(n) + br(n-T)$
- zpětná operace: $r(n) = e(n) - br(n-T)$
- říká se mu „dlouhodobý“, protože prediktuje ze vzorků vzdálených až 20ms, zatímco krátkodobý prediktor (LP) prediktuje ze vzorků vzdálených ~ 2 ms
- gain b je záporný

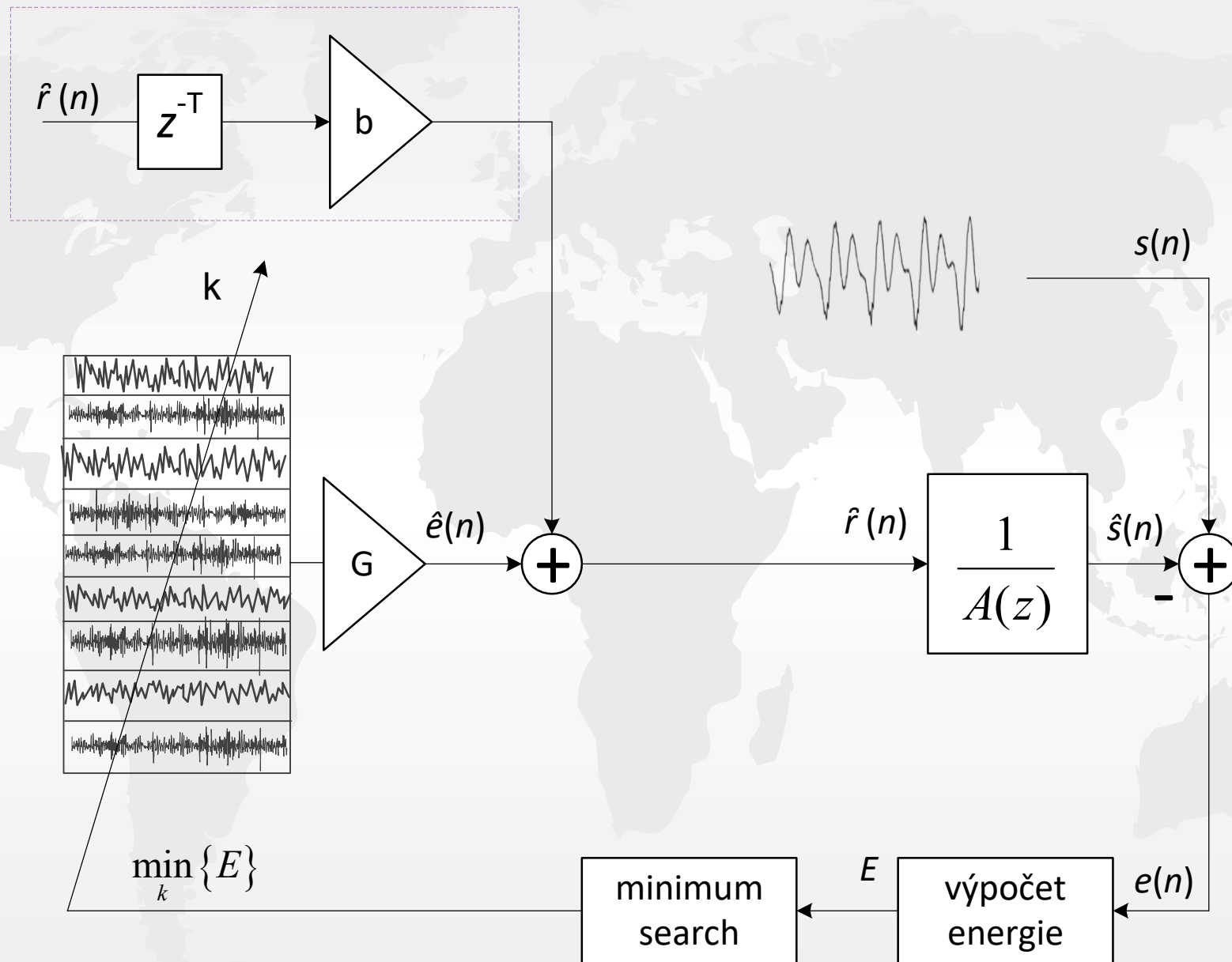
Dlouhodobý prediktor (LTP)



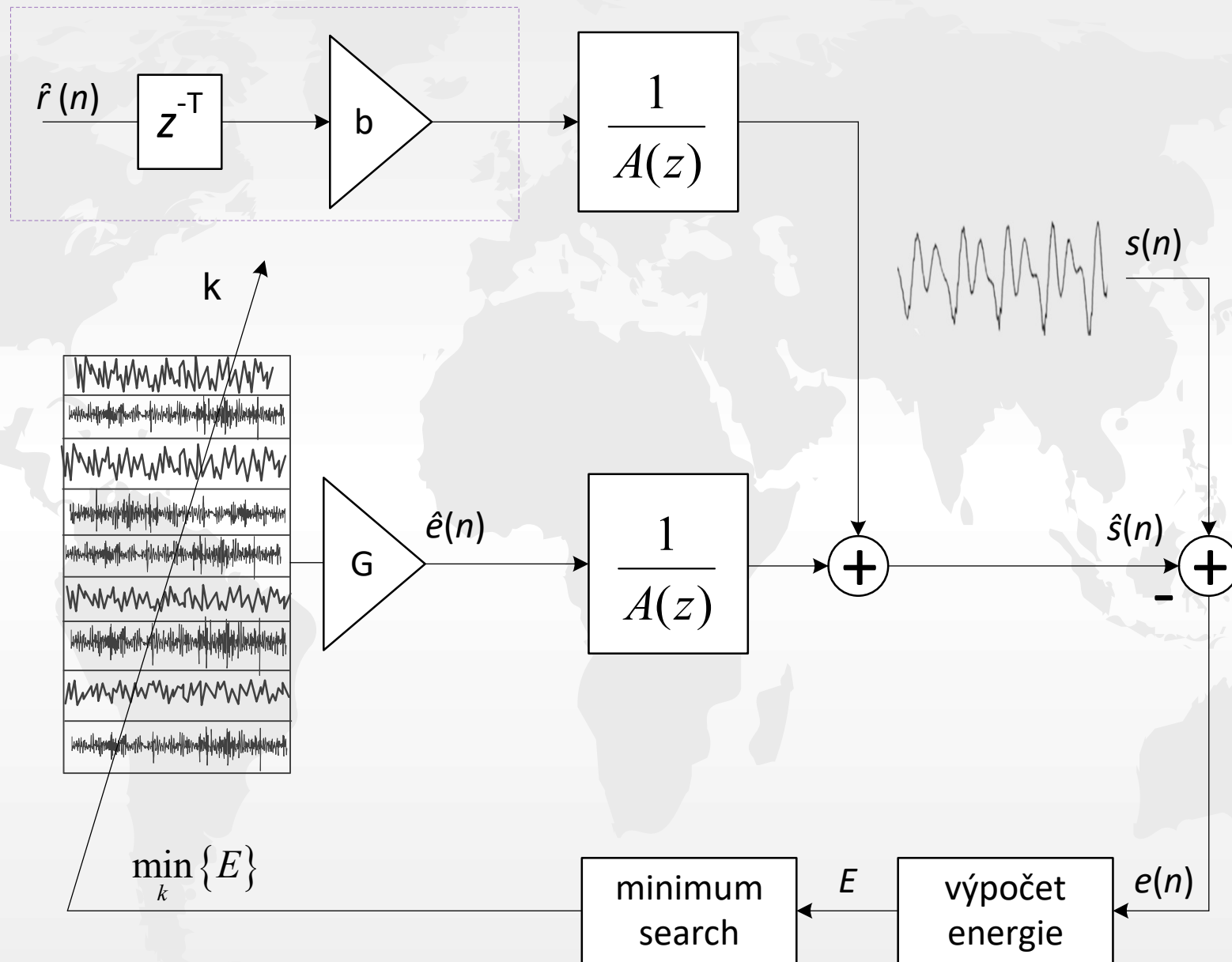
Dlouhodobý prediktor (LTP)



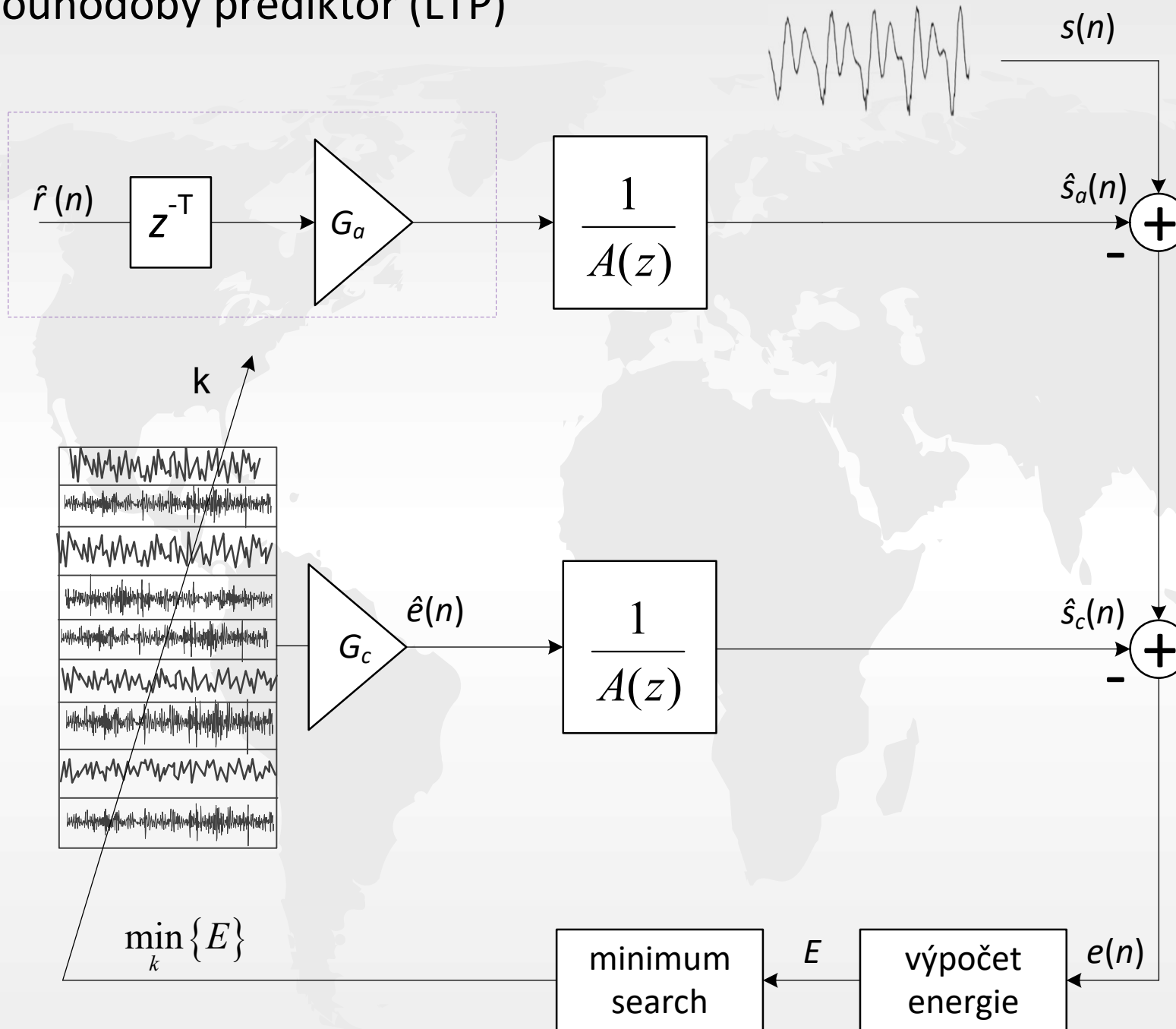
Dlouhodobý prediktor (LTP)



Dlouhodobý prediktor (LTP)

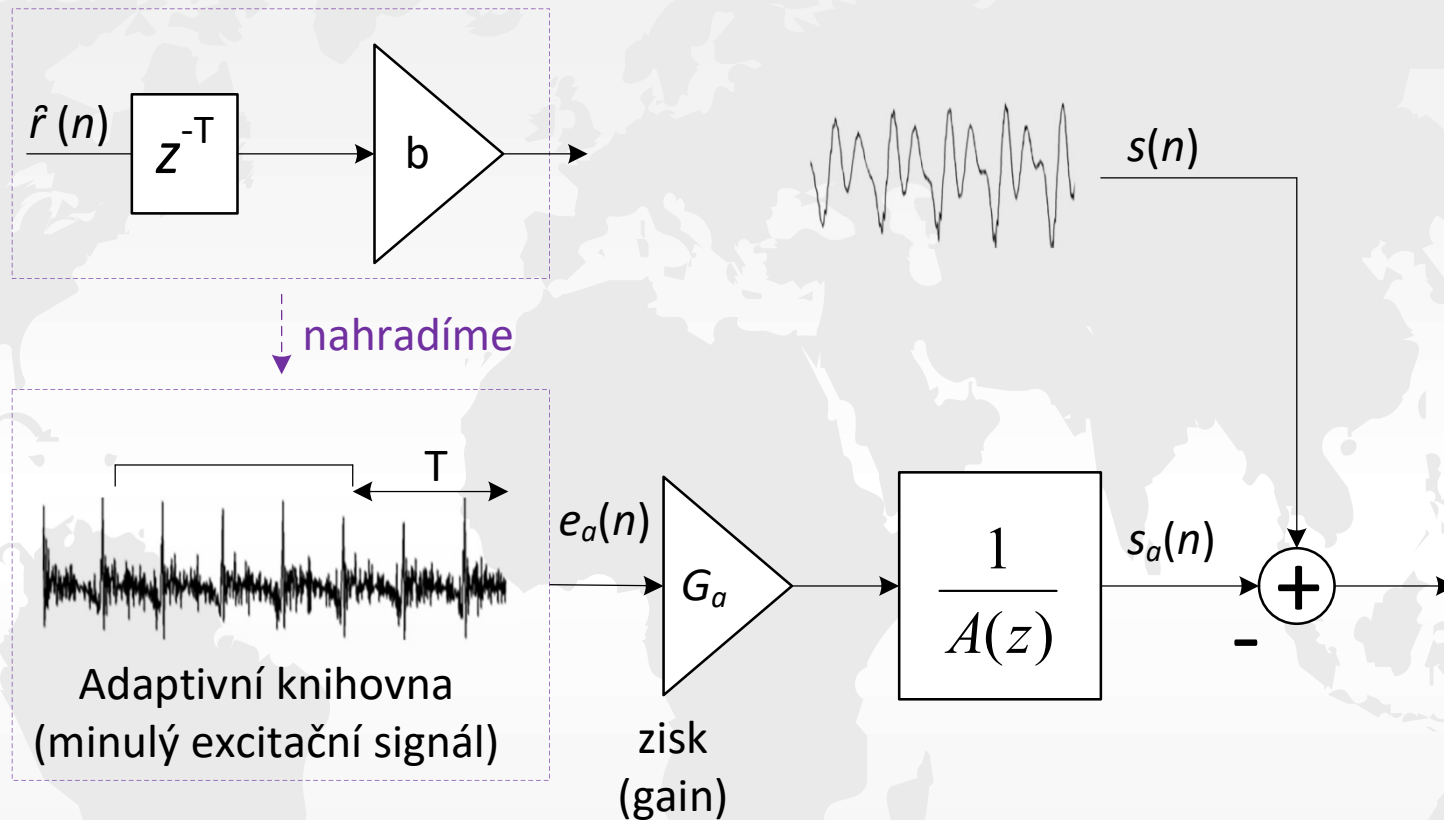


Dlouhodobý prediktor (LTP)



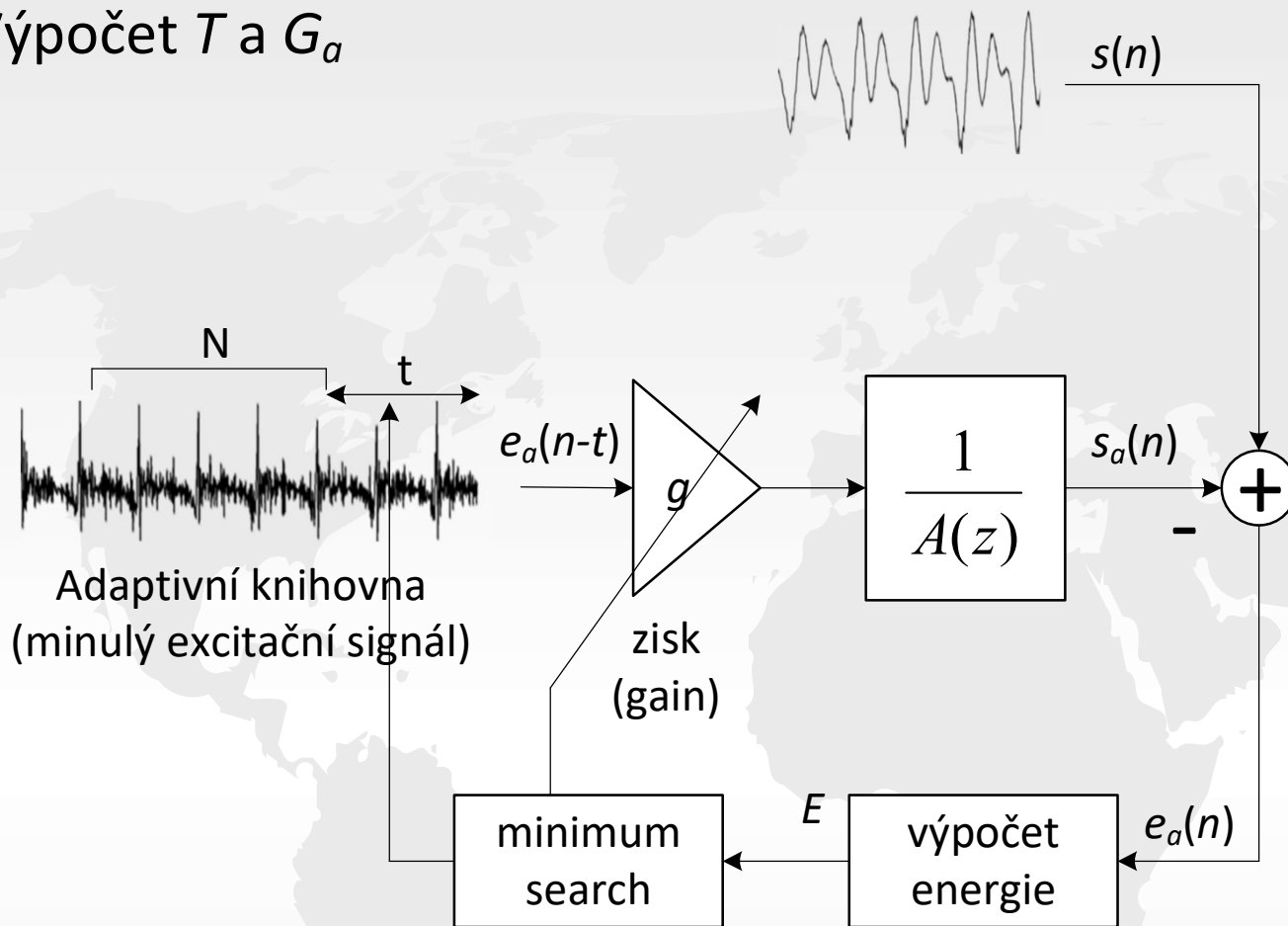
- přejmenování
 $b \rightarrow G_a$
 $G \rightarrow G_c$
 $\hat{s}(n) \rightarrow \hat{s}_c(n)$

Zavedení adaptivní knihovny



- dlouhodobý prediktor lze nahradit tzv. adaptivní knihovnou, což je v podstatě minulý excitační signál (pozn. $r(n)$ je LP reziduum, $\hat{r}(n)$ je excitace)

Výpočet T a G_a



postup:

- zvolíme t a vybereme vektor $e_a(n-t)$
- vypočítáme g

$$g = \frac{-\sum_{n=1}^N r(n)e_a(n-t)}{\sum_{n=1}^N e_a^2(n-t)}$$

- vypočítáme $s_a(n)$

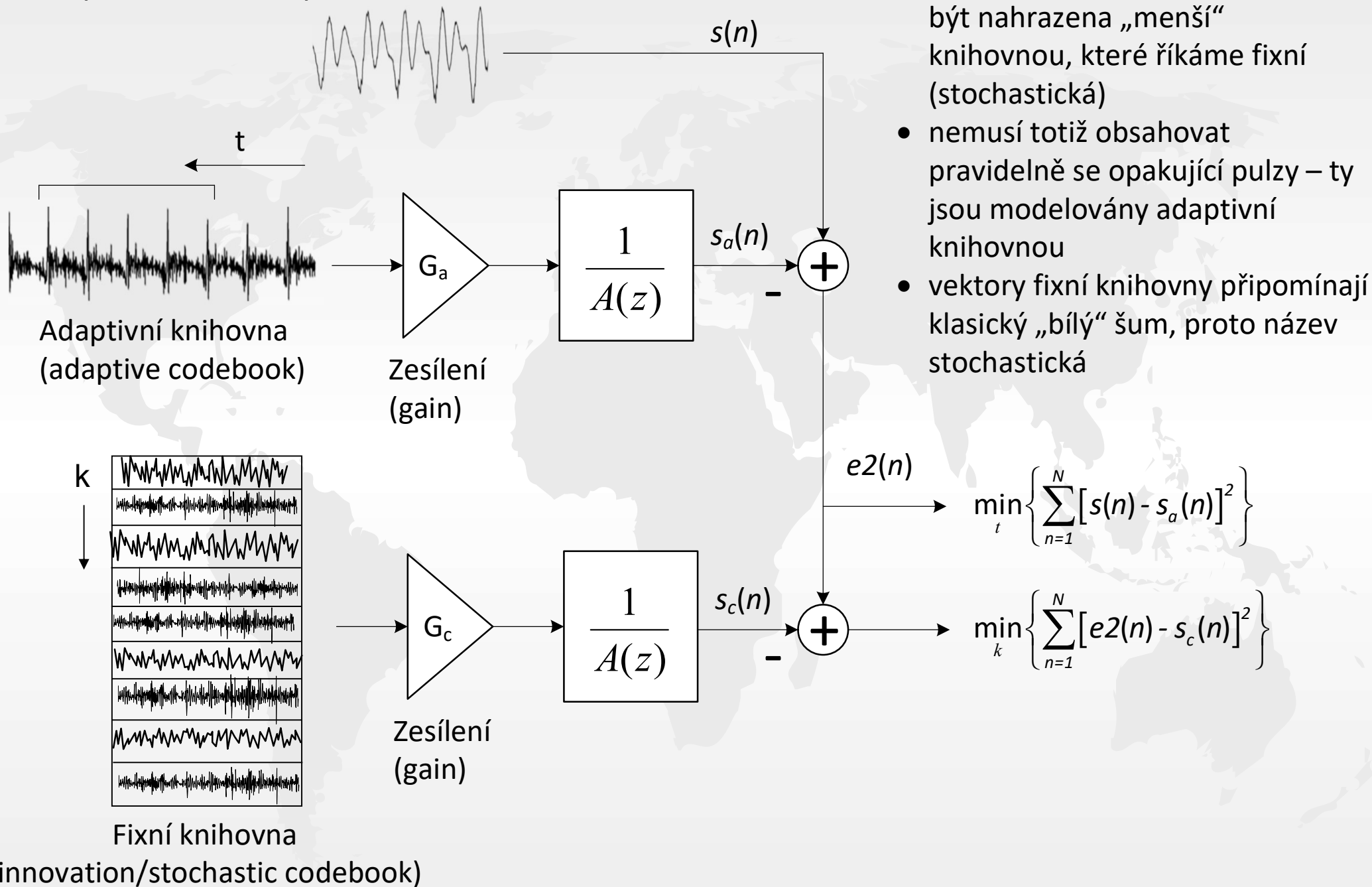
$$s_a(n) = ge_a(n-t) - \sum_{m=1}^M a_m s_a(n-m)$$

- vypočítáme E

$$E = \sum_{n=1}^N [s(n) - s_a(n)]^2$$

- parametry T a G_a najdeme tak, abychom minimalizovali chybu mezi $s(n)$ a $s_a(n)$
- iterativně „prohledáváme“ adaptivní knihovnu a testujeme N -dlouhé úseky minulé excitace různě vzdálené od počátku ($n=0$)
- abychom ušetřili počet nutných iterací, prohledáváme pouze v okolí základního tónu
- pro každé testované t dopočítáme analyticky optimální hodnotu gainu g
- gain g se nekóduje, později ho totiž přepočítáme

Fixní (stochastická) knihovna



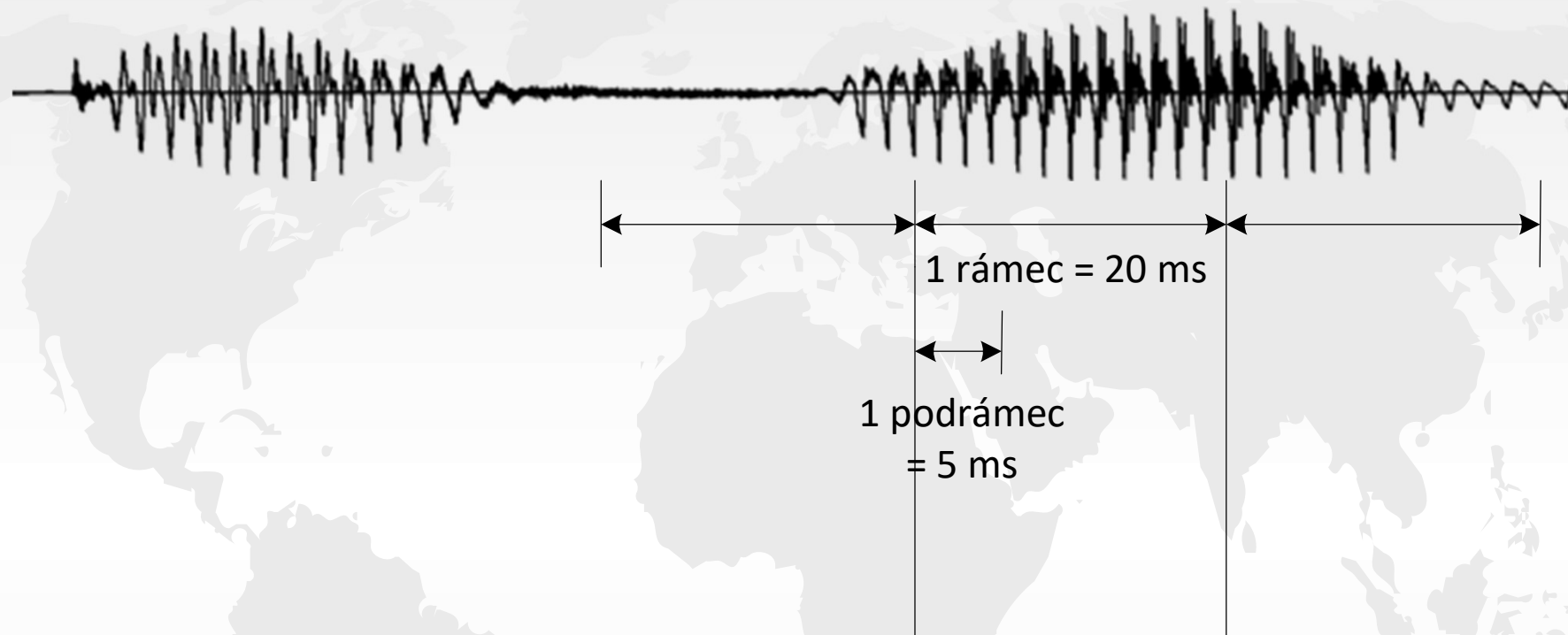
- původní „obří“ knihovna může být nahrazena „menší“ knihovnou, které říkáme fixní (stochastická)
- nemusí totiž obsahovat pravidelně se opakující pulzy – ty jsou modelovány adaptivní knihovnou
- vektory fixní knihovny připomínají klasický „bílý“ šum, proto název stochastická



Rámce a podrámce

Rámce a podrámce

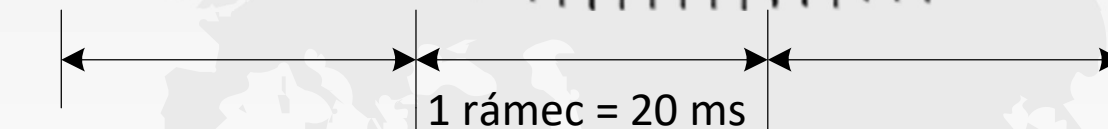
Řečový signál



- klasické rámce o délce 15-25ms jsou pro některé parametry (LP koeficienty, T_0) příliš dlouhé
- lidé totiž mění polohu vokálního traktu nebo základní tón někdy mnohem rychleji
- řešením je tedy zkrácení rámce např. na 5ms (podrámec) pro výpočet těchto parametrů
- ale pozor, např. LP analýza je docela komplexní a dělat ji 4x v jednom rámci stojí hodně MIPS

interpolace LP koeficientů v podrámcích

Řečový signál



1 rámeček = 20 ms

1 podrámeček = 5 ms

Hamming window
pro LP analýzu

asymetrické okno
pro LP analýzu

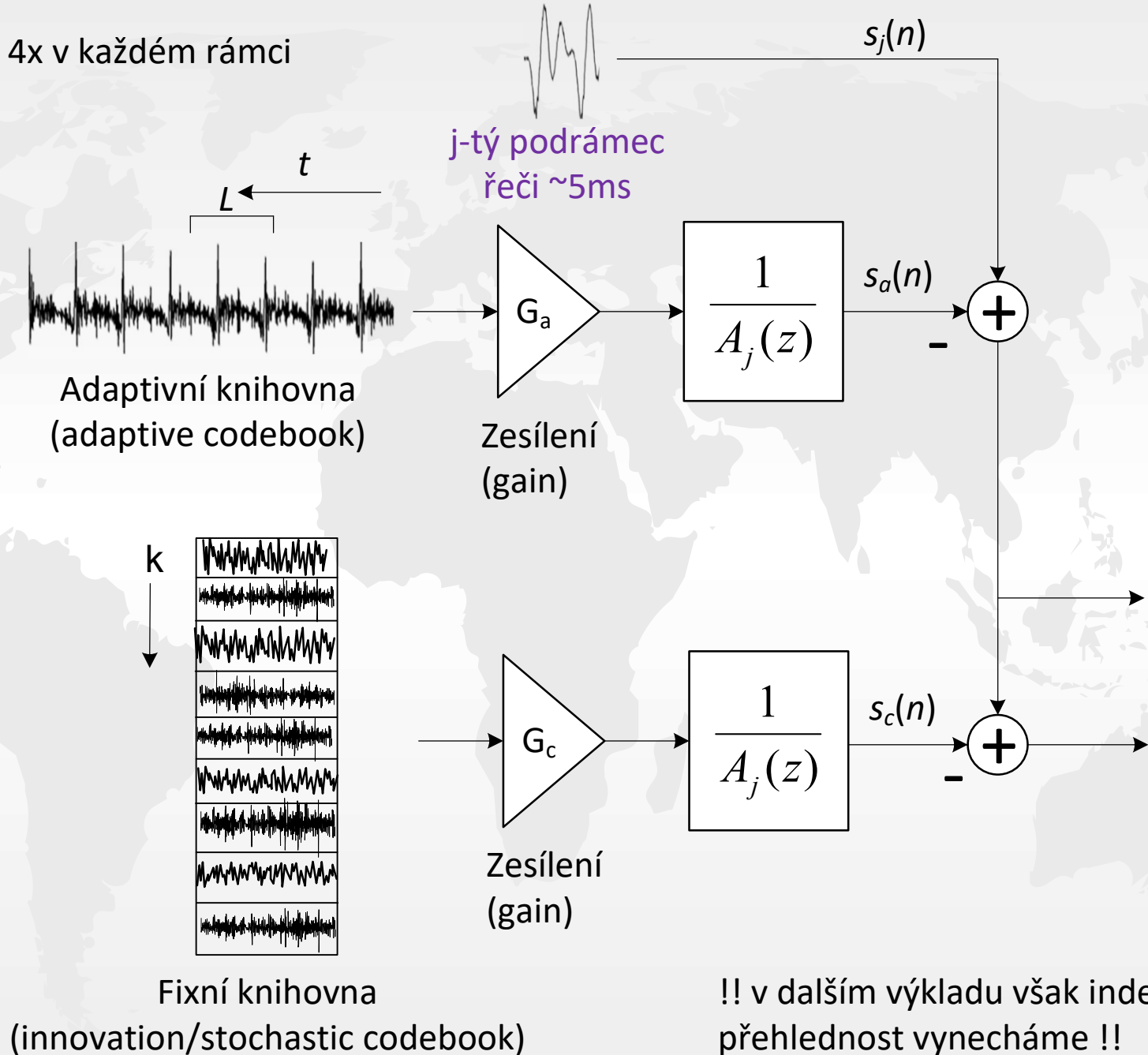


$$0.8 * a_i + 0.2 * \text{old}_a_i$$

- lepší nápad – budeme dělat LP analýzu jen jednou, např. na konci rámeče a 3x interpolovat mezi starými a novými koeficienty LP filtru
- k tomu ale potřebujeme jiný tvar okna -> asymetrické okno
- asymetrické okno musí co nejméně zasahovat do „budoucího“ rámeče, protože kodek musí na tyto vzorky „čekat“
- tyto vzorky „budoucího“ rámeče nazýváme LOOKAHEAD (způsobuje zpoždění kodeku, tzv. algorithmic delay)

CELP v podrámcích

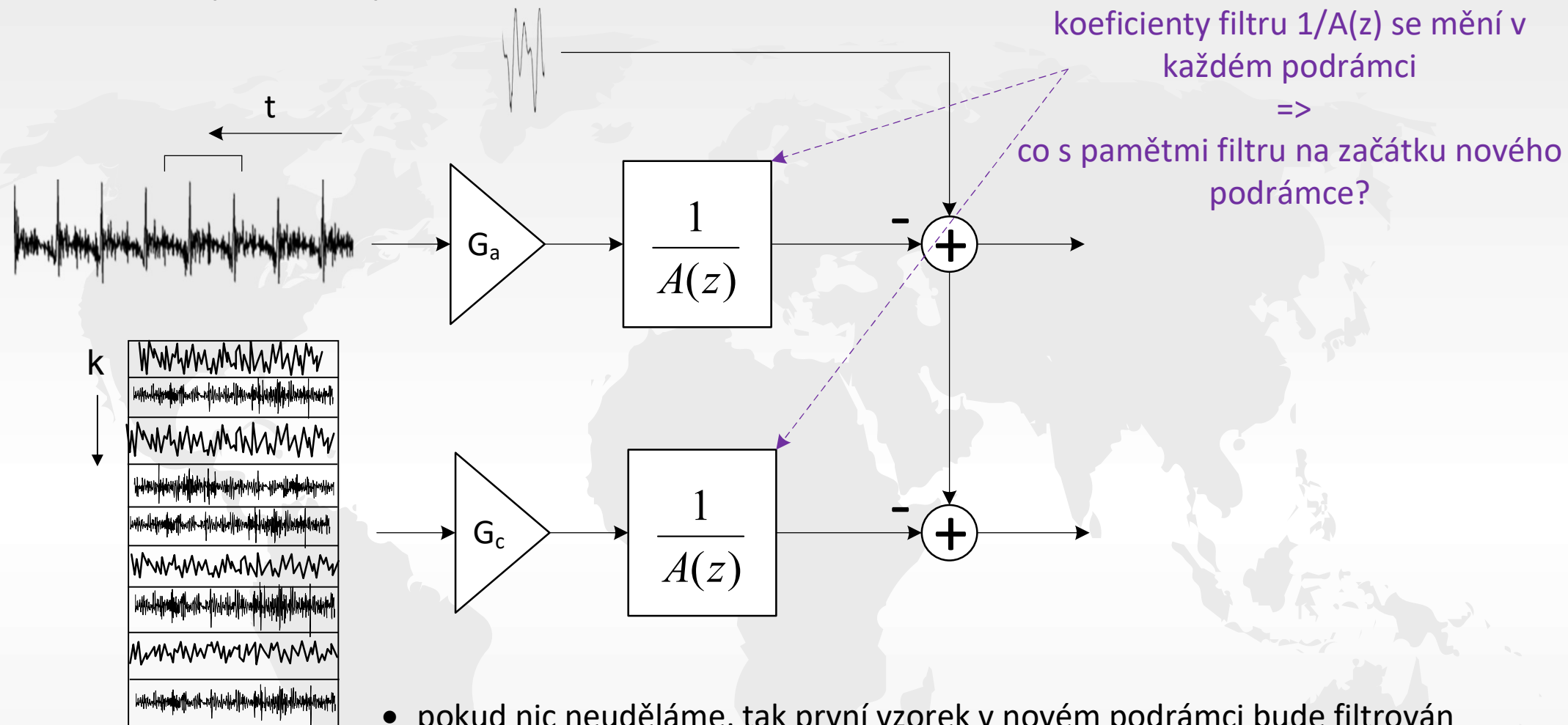
tohle musíme dělat 4x v každém rámcí



A faint, light gray world map is visible in the background of the slide, centered behind the text.

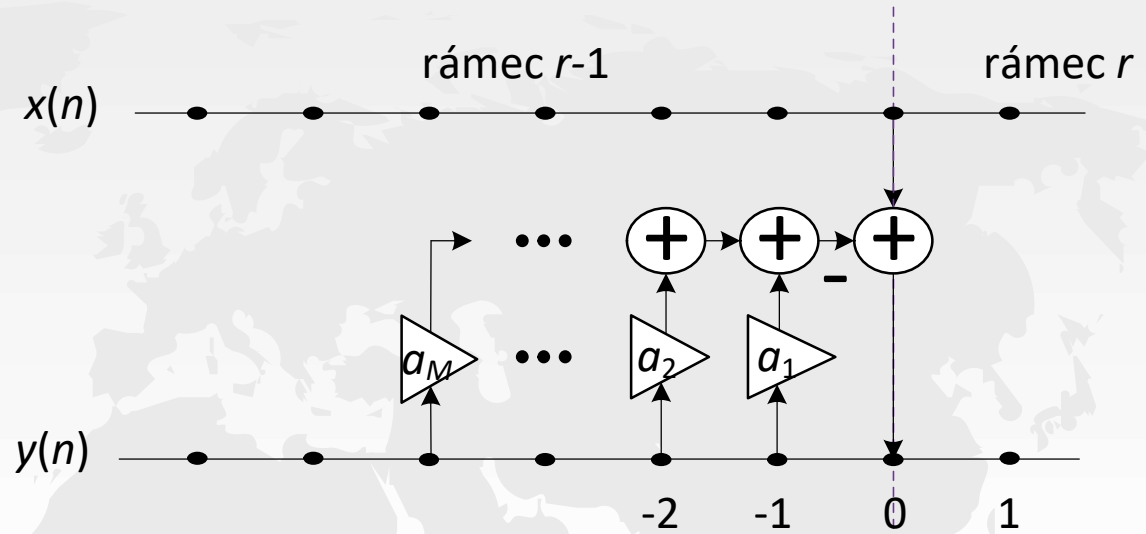
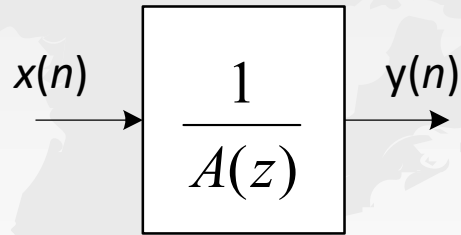
Zero-Input Response Zero-State Response

Problém s přechody mezi rámcí



- pokud nic neuděláme, tak první vzorek v novém podrámci bude filtrován nesprávně v důsledku špatných pamětí a v syntetickém signálu vznikne skok (discontinuity)
- zkusíme rozdělit odezvu filtru na dvě části – jedna, která bude odpovídat pouze na minulý výstupní signál (uložen v pamětech filtru) a druhá, která bude odpovídat pouze na vstupní signál
- obě složky sečteme

Zero-Input Response (ZIR) a Zero-State Response (ZSR)



$$y(n) = x(n) - \sum_{i=1}^M a_i y(n-i) \quad 0 \leq n < N$$

pro $M=1$ máme:

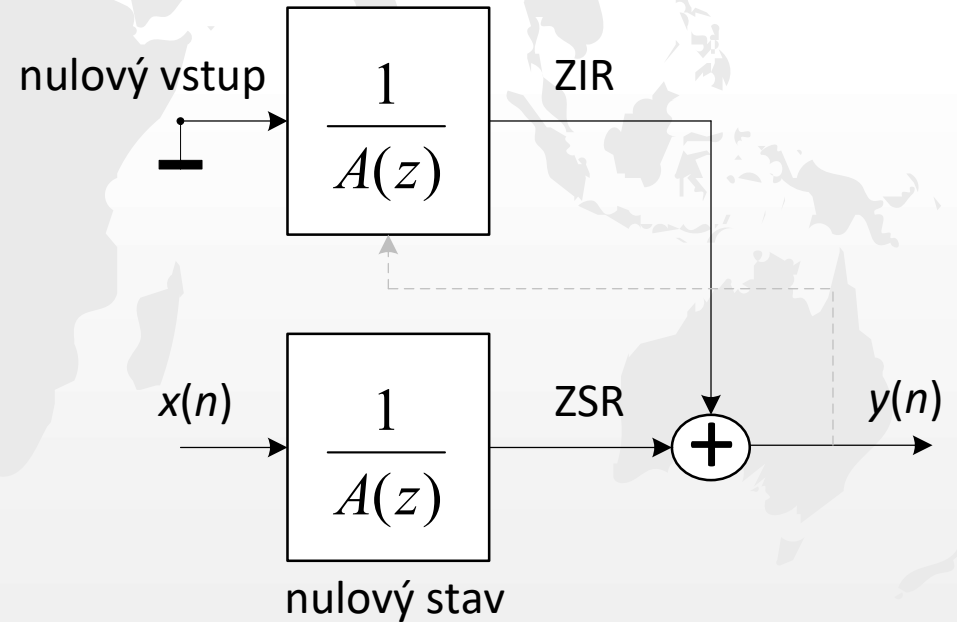
$$y(0) = x(0) - a_1 y(-1)$$

$$y(1) = x(1) - a_1 x(0) + a_1^2 y(-1)$$

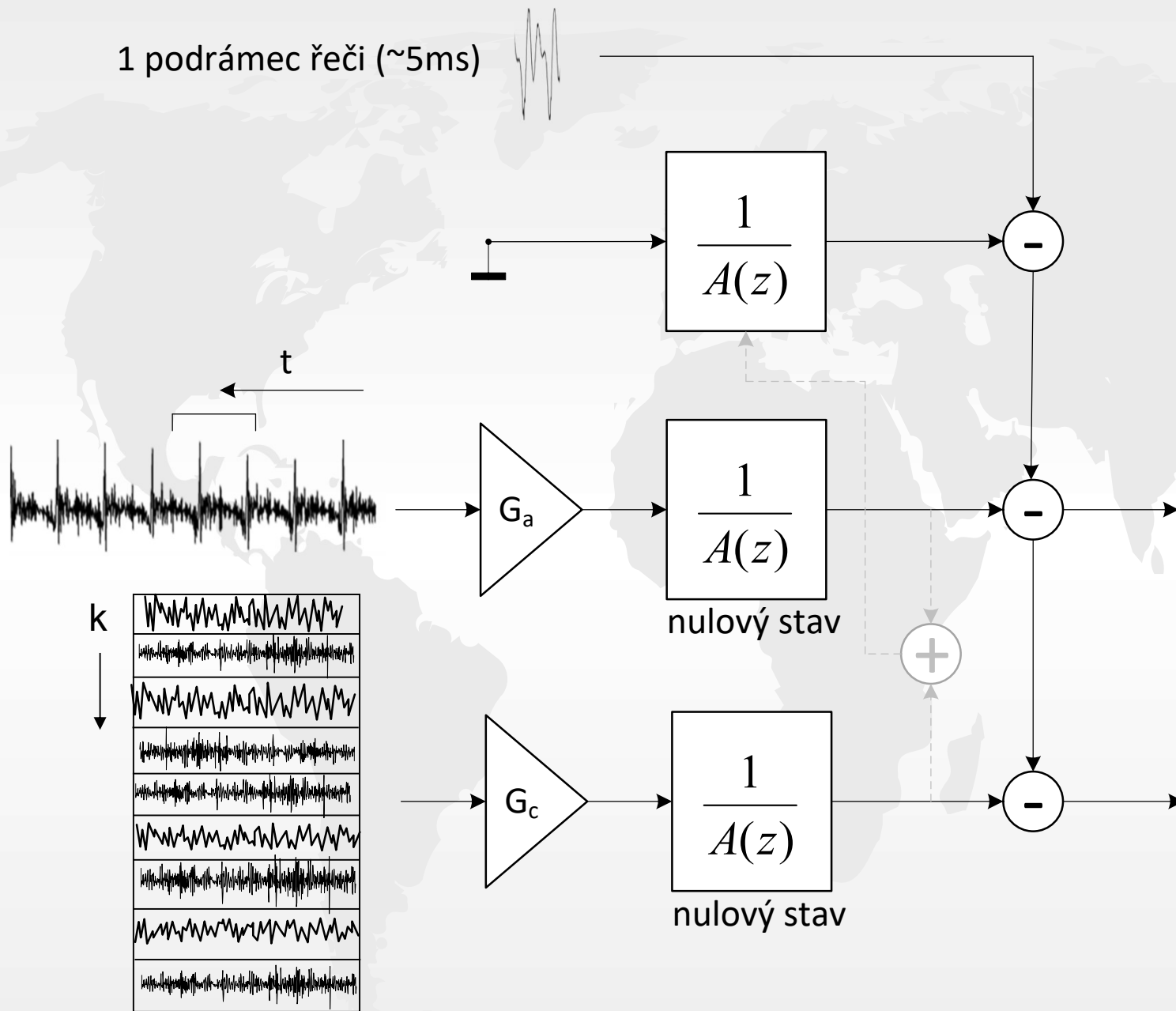
$$y(2) = x(2) - a_1 x(1) + a_1^2 x(0) - a_1^3 y(-1)$$

ZSR

ZIR



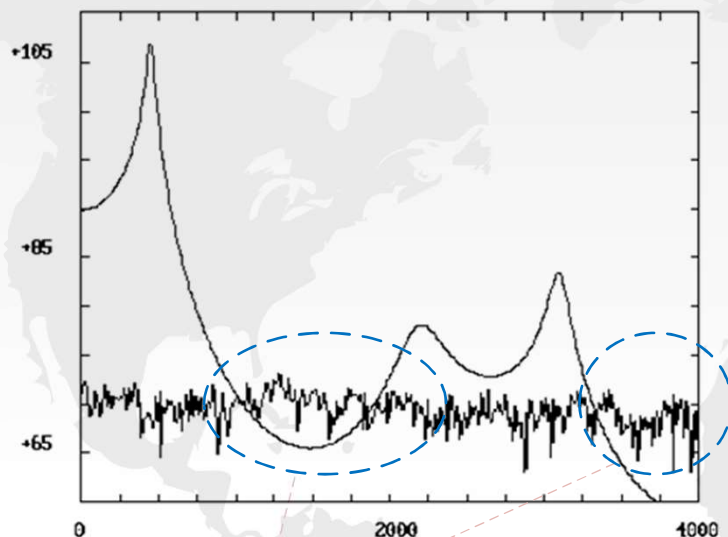
Odečtení ZIR



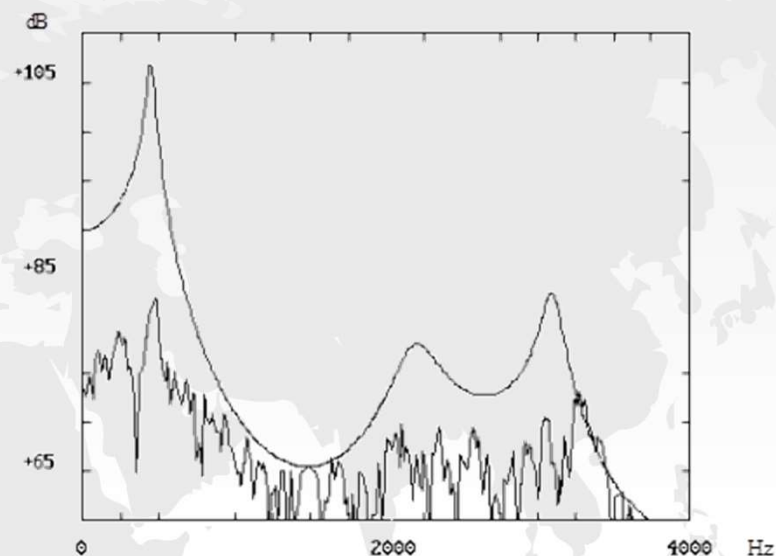


Perceptuální váhování

Maskování kvantizačního šumu



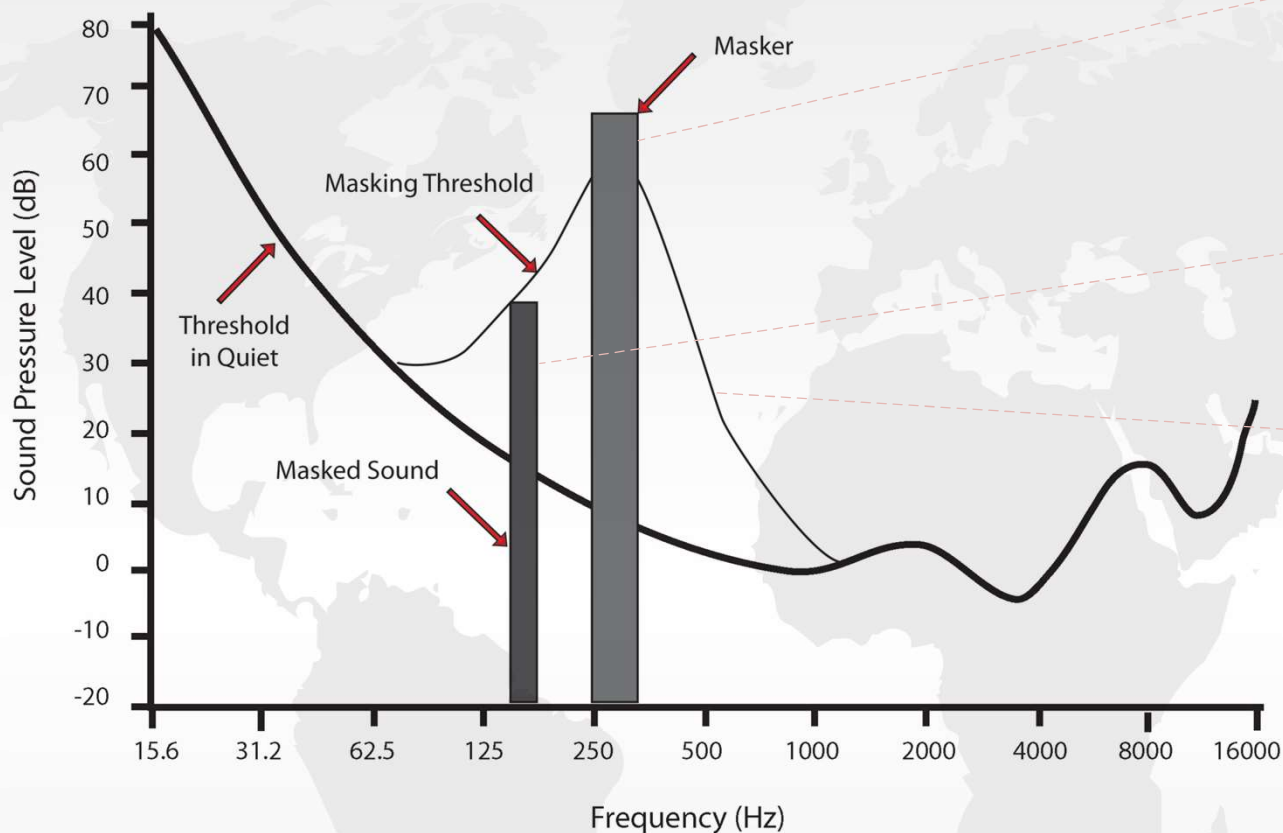
tohle bude slyšet



chtěli bychom, aby kvantizační šum vypadal nějak takhle

- Lidské ucho nedokáže registrovat zvuky, které jsou maskovány silnějším signálem
- i když najdeme adaptivní a fixní část excitace jak nejlépe umíme, pořád zbyde nějaký chybový signál, tzv. kvantizační šum
- části kvantizačního šumu v „údolích“ mohou být slyšet
- potřebovali bychom tenhle šum „zdeformovat“ tak, aby potlačen v „údolích“ a na oplátku si můžeme dovolit, aby byl zesílen v oblasti formantů

Maskování tónů



maskovací šum o centrální frekvenci 300Hz a šířce pásma 100Hz a úrovni 65 dB

maskovaný signál o centrální frekvenci 200Hz a šířce pásma 50Hz a úrovni 38 dB

práh slyšitelnosti

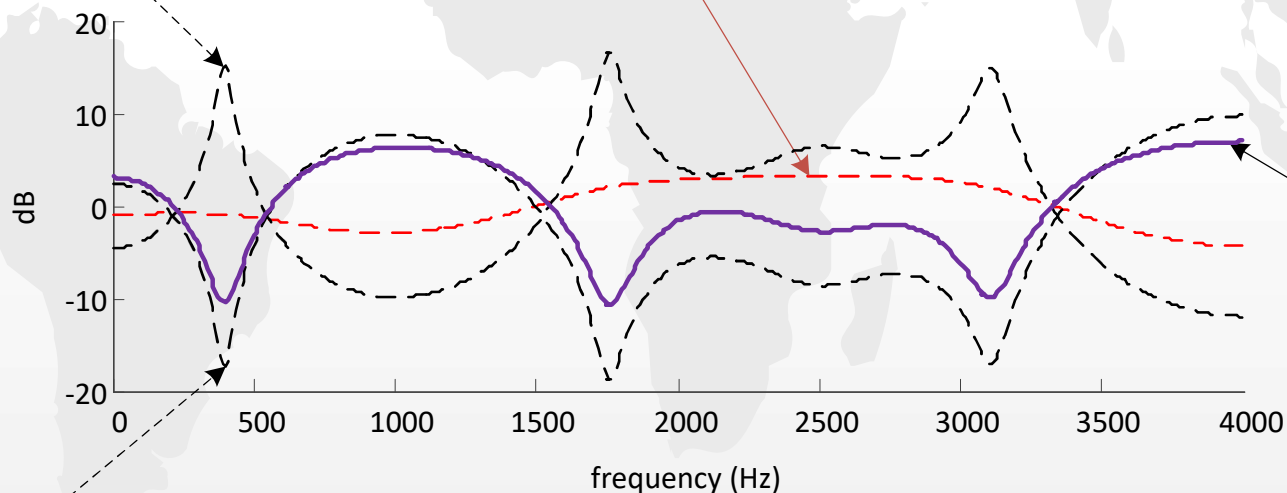
- tóny v blízkosti silného tónu jsou maskovány
- spektrální komponenty s úrovní pod prahem slyšitelnosti není třeba kódovat
- lze tolerovat vyšší úroveň kvantizačního šumu v blízkosti silných tónů, např. formantů
- demo na <https://www.youtube.com/watch?v=k6DVywW5NR4>

Perceptuální váhování

$$W(z) = \frac{A(z)}{A(z/\gamma)} = \left[1 - \sum_{i=1}^M a_i z^{-i} \right] \left[\frac{1}{1 - \sum_{i=1}^M \gamma^i a_i z^{-i}} \right]$$

$$\frac{1}{A(z)} = \frac{1}{1 - \sum_{i=1}^M a_i z^{-i}} \quad \frac{1}{A(z/\gamma)} = \frac{1}{1 - \sum_{i=1}^M \gamma^i a_i z^{-i}}$$

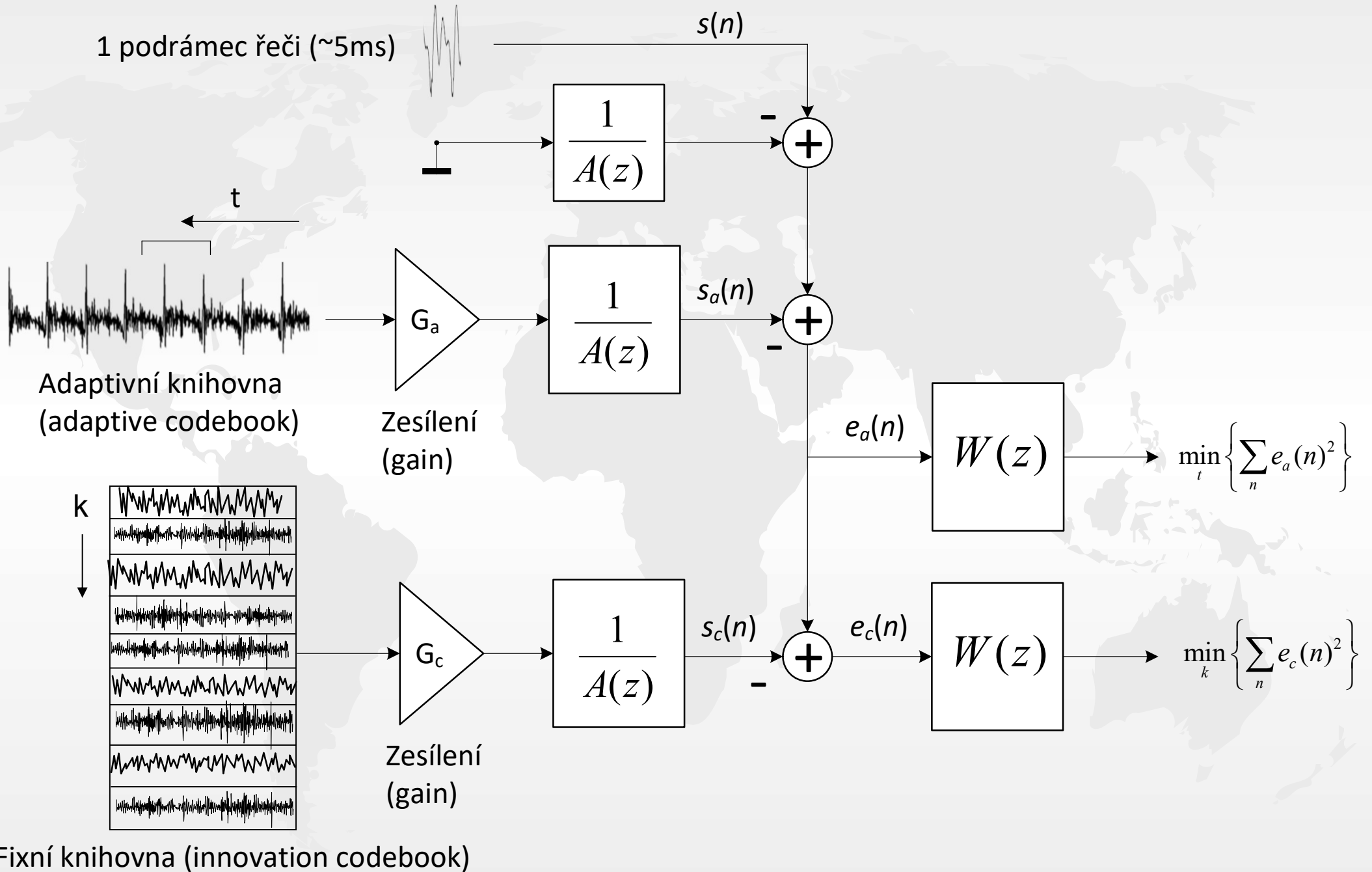
- musíme najít takový tvar perceptuálního (váhovacího) filtru, který bude dávat velkou váhu „údolím“ a malou váhu „formantům“
- celkový zisk filtru musí být 1, jinak by docházelo k postupnému zesilování/ zeslabování signálu
- $\gamma = 0.8 - 0.9$



$W(z)$
váhovací filtr

$$A(z) = 1 - \sum_{i=1}^M a_i z^{-i}$$

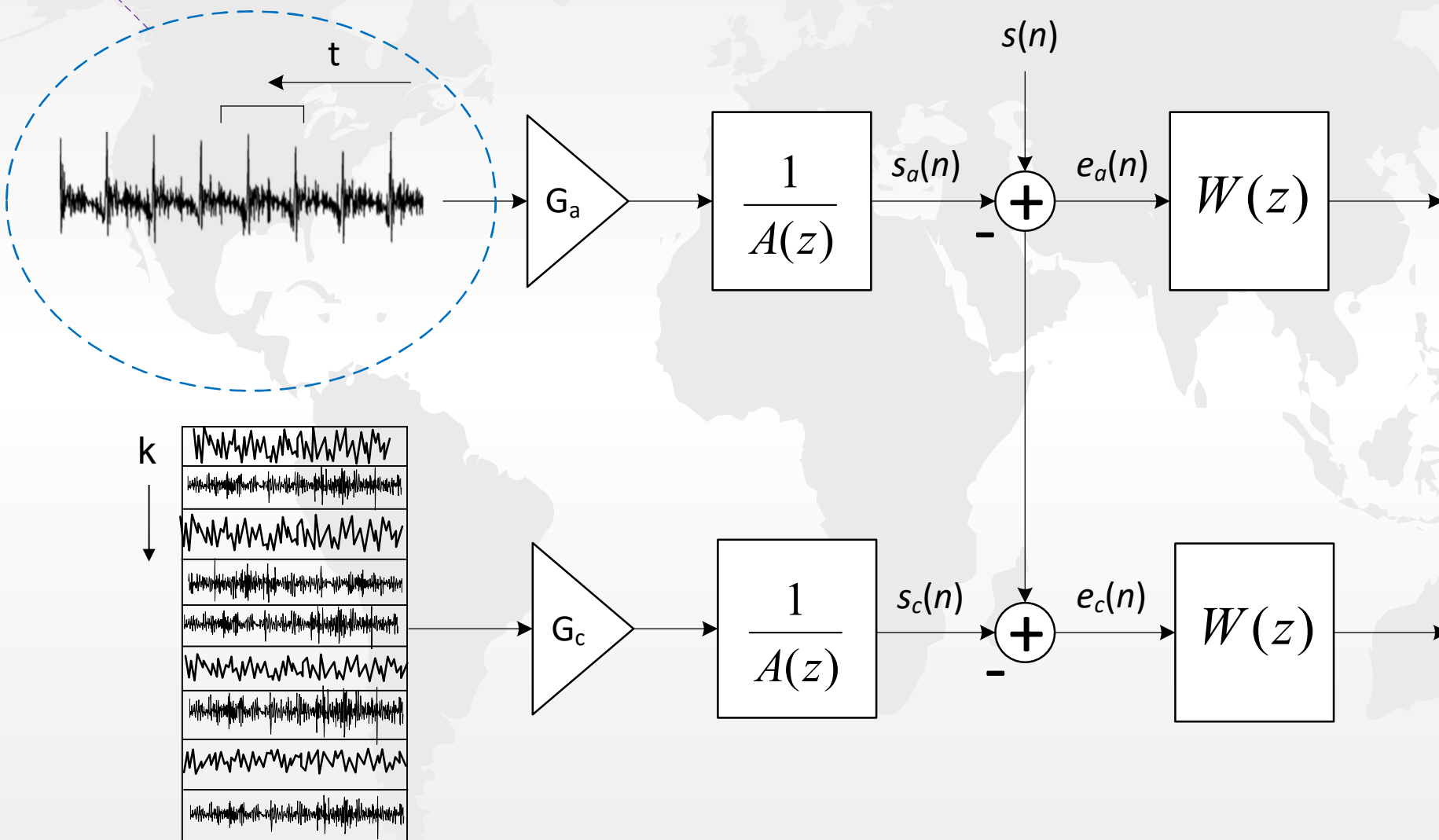
Zavedení perceptuálního filtru



Problém s výpočetní náročností

Adaptivní knihovna: (7-9 bitů, t.j. 128 – 512 vektorů)

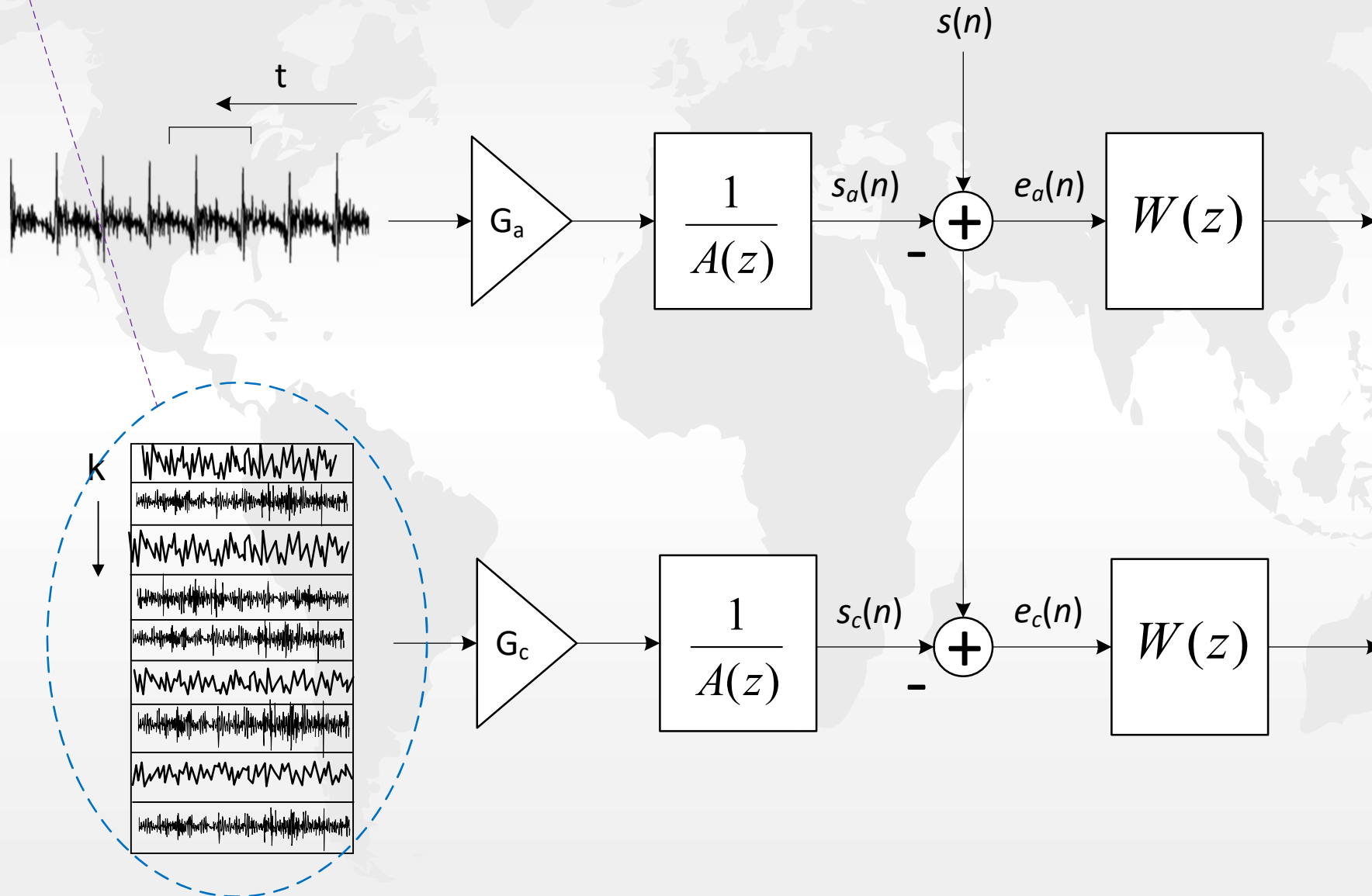
(prohledáváme knihovnu pouze kolem základního tónu)



Problém s výpočetní náročností

Fixní knihovna: (10-88 bitů, t.j. $1024 - 3 \cdot 10^{26}$ vektorů)

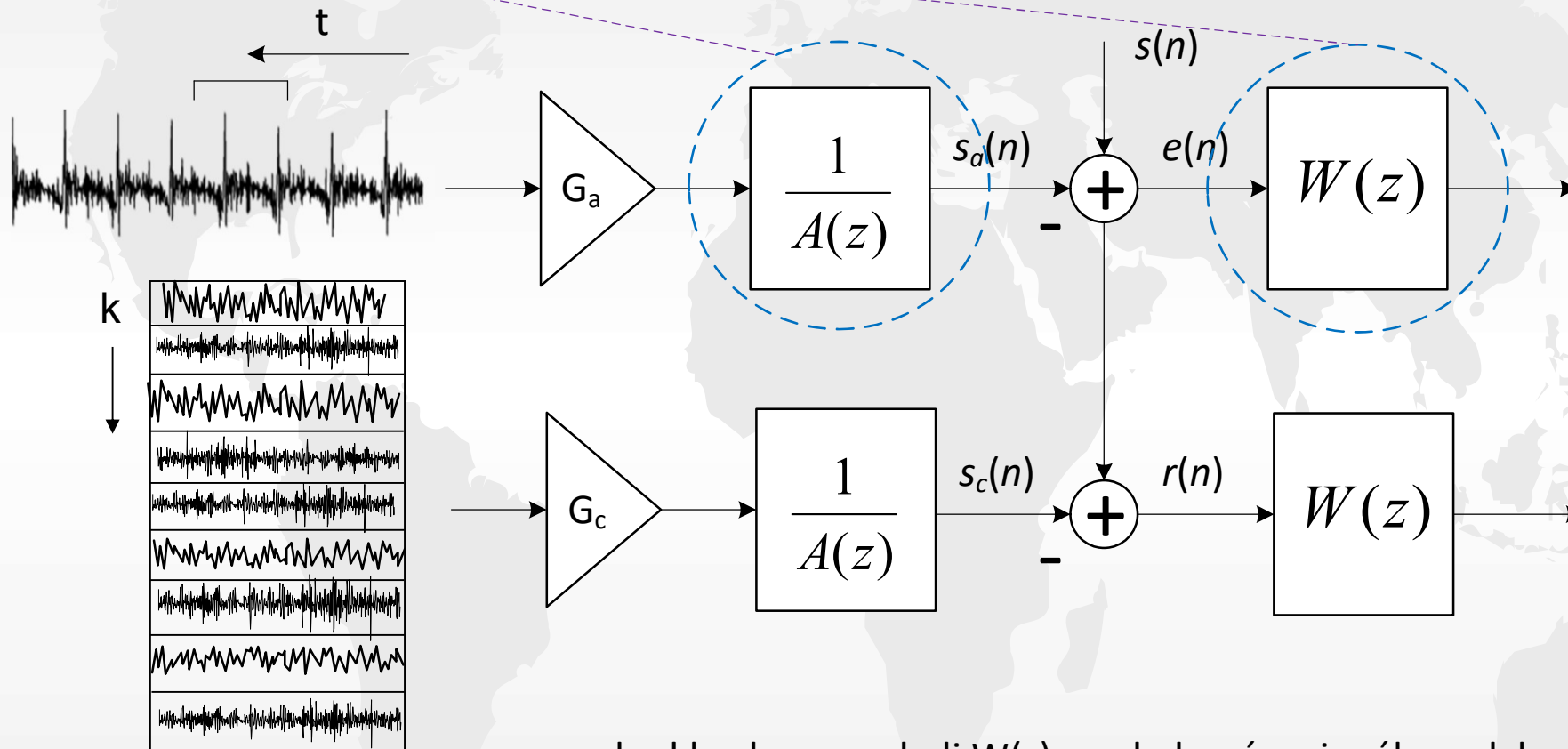
(vnucení jednoduchých struktur – pouze několik pulzů na stopu, omezený počet pozic pulzů, znaménka)



Problém s výpočetní náročností

Filtrace: (pro každý codevector nutno provést konvoluci s filtry 16.řádu)

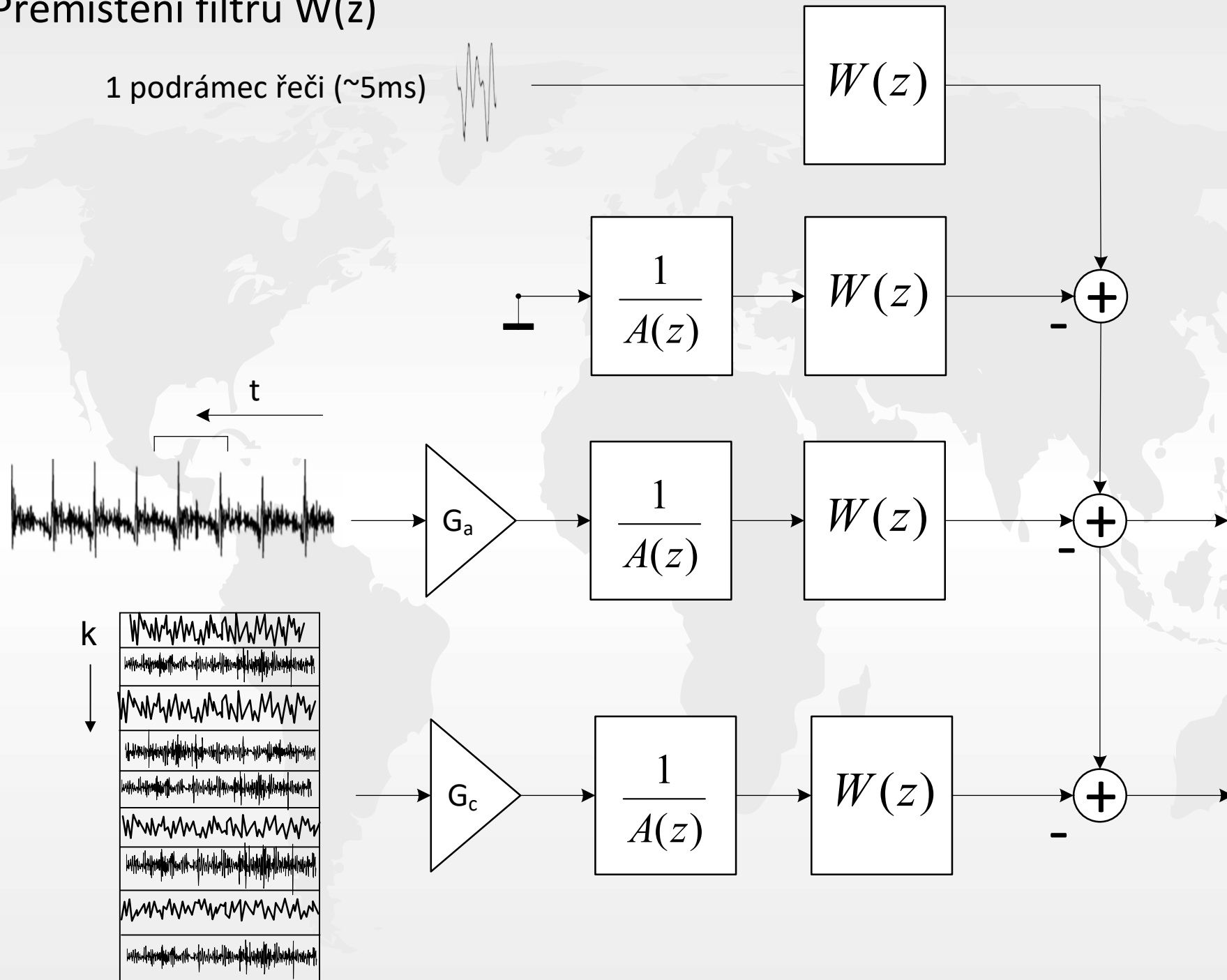
(nahrazení filtrů $1/A(z)$ a $W(z)$ jejich impulzní odezvou) – viz dále



- pokud bychom nechali $W(z)$ na chybovém signálu, pak bychom museli počítat konvoluci pro každý testovaný vektor z knihoven -> to je moc výpočetně náročné
- lepší nápad je přesunout filtr $W(z)$ do obou větví excitace a na vstup

Přemístění filtru $W(z)$

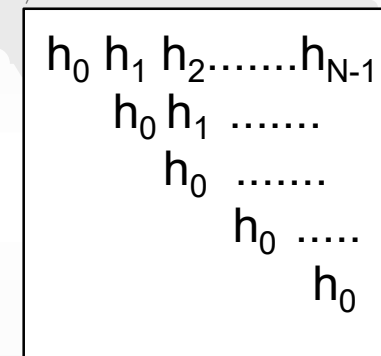
1 podrámec řeči (~5ms)



Nahrazení filtrů impulzní odezvou

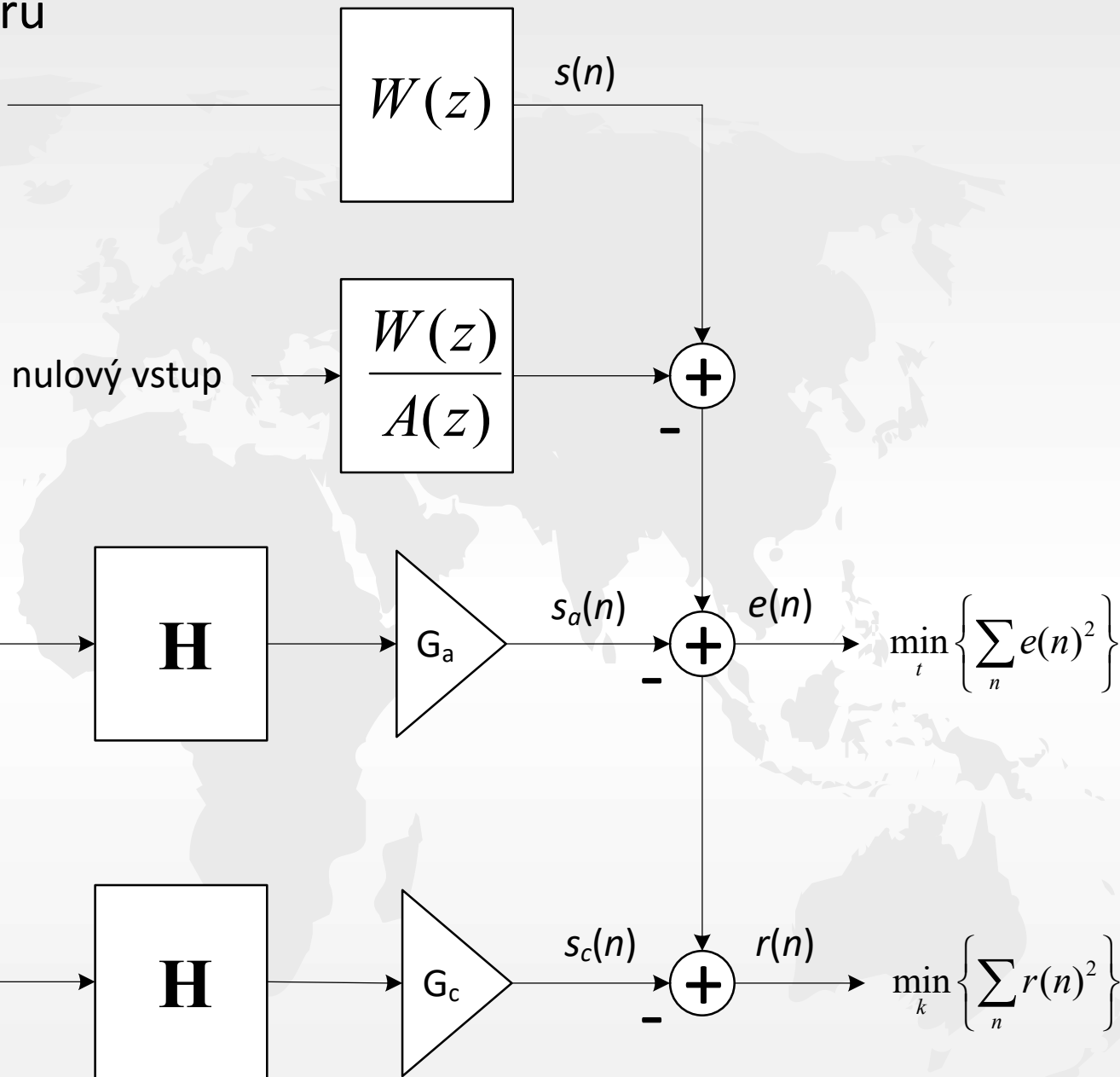
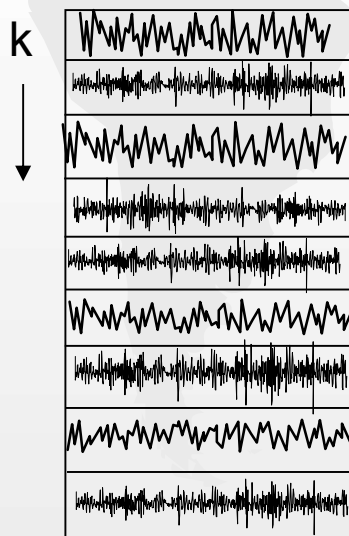
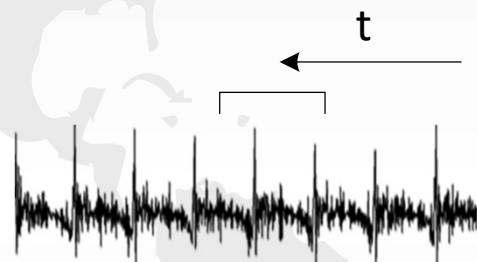
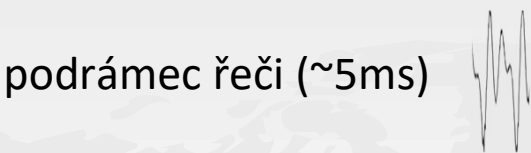


- kaskádu dvou filtrů nahradíme obyčejným maticovým násobením s impulzní odezvou filtru $W(z)/A(z)$
- jenže impulzní odezva IIR filtru je nekonečně dlouhá, tak nám nezbyvá nic jiného než ji „ustříhnout“ na konci rámce
- tím vznikne chyba, která ovšem se vzdáleností od začátku rámce klesá
- koeficienty h_0, h_1, \dots, h_{N-1} tvoří impulzní odezvu filtru $W(z)/A(z)$
- vzhledem k předpokladu nulového stavu pamětí má matice H triangulární tvar

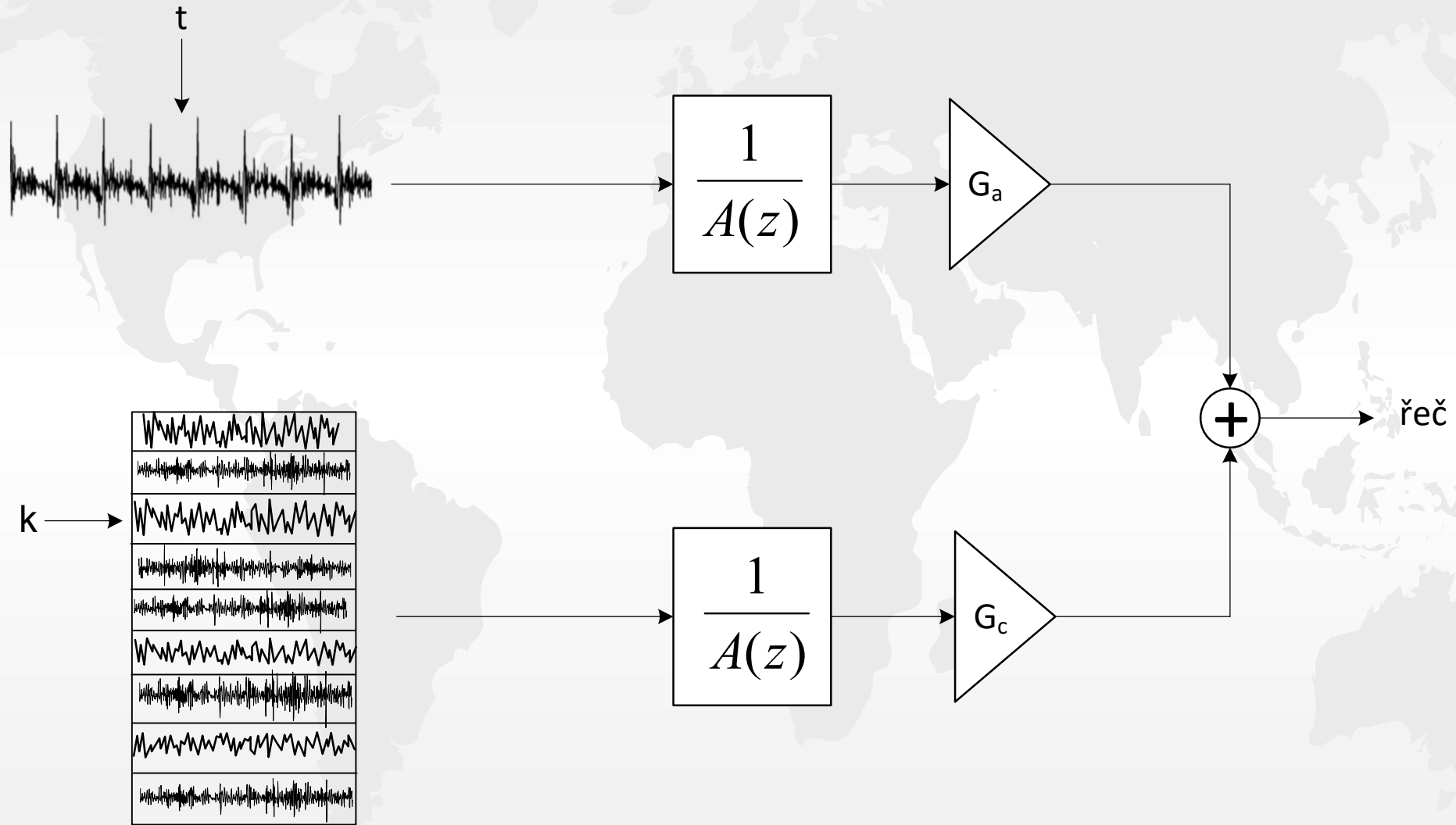


Celkové schéma CELP enkodéru

1 podrámec řeči (~5ms)



Celkové schéma CELP dekodéru



Kodeky používající CELP

CELP

4.8 kbps, 19 MIPS, U.S. DoD (Dept. of Defence) FS 1016 (1991)

lin. knihovna s hodnotami 0, 1, -1

jednotlivé vzorky se liší pouze ve 2 vzorcích

RPE-LTP – Regular Pulse-Excited Long-Term Prediction

13 kbps, 18.5 WMOPS, ETSI GSM FR (1989)

buzení podvzorkováno faktorem 3 a je kvantována pouze poloha prvního impulzu
další impulzy jsou kvantovány pomocí APCM

VSELP – Vector Sum Excitation LP

7.95 kbps, 20 MIPS, TIA IS-54 (1989)

5.6 kbps, 30 MIPS, ETSI GSM HR (1995)

několik vektorů, filtrovaných předem, tvoří bázi

výsledné vektory se pak vytvoří jejich lineární kombinací



ACELP

ACELP

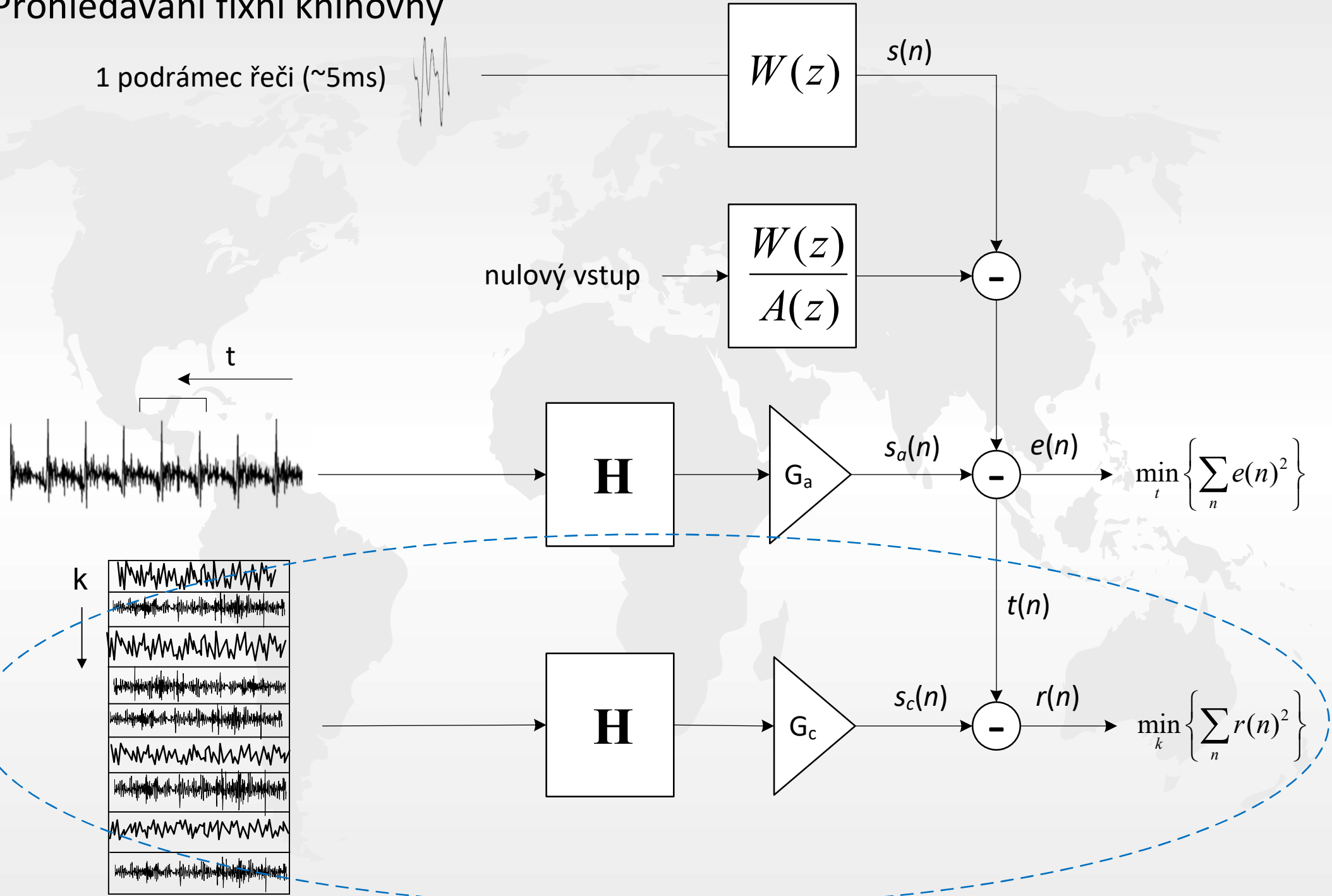
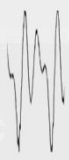
- ACELP® je patentovanou technologií VoiceAge Corp. a Université de Sherbrooke, CANADA
- vyvinuto výzk. skupinou Jean-Pierre Adoula v roce 1987
- poprvé publikováno zde: Adoul, J-P. and C. Lamblin (1987). “A Comparison of Some Algebraic Structures for CELP Coding of Speech,” IEEE ICASSP, pp. 1953–1956.



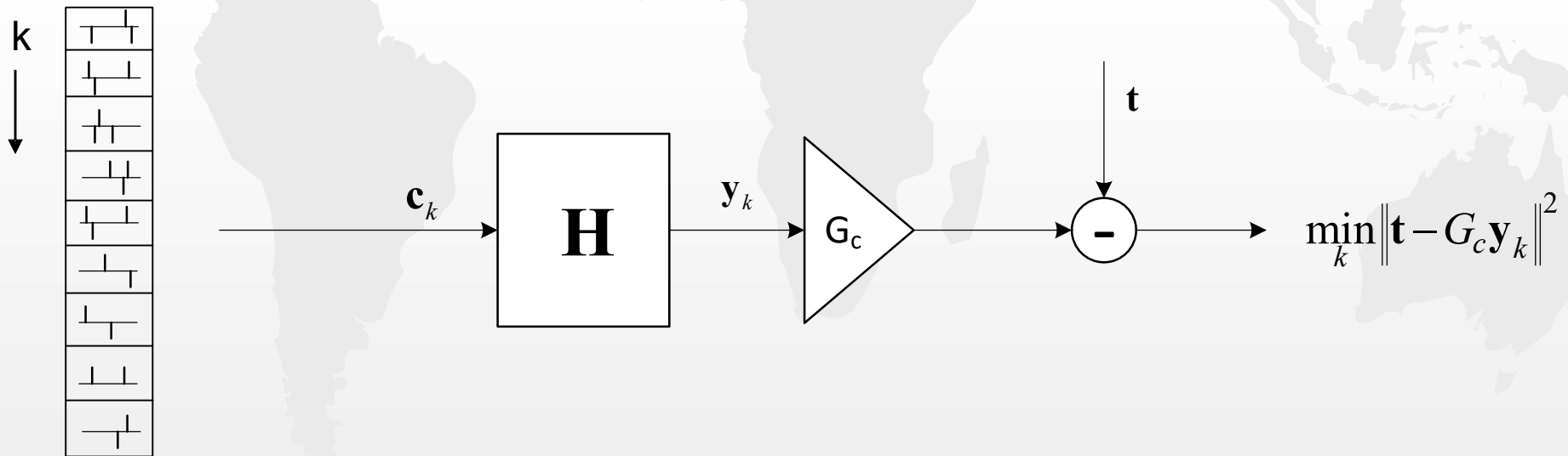
- kouzlo ACELPu spočívá v tom, že dokáže nahradit „obří“ fixní knihovnu signálů jednoduchou knihovnou s algebraickou strukturou, kde je jen několik málo pulzů v přesně definovaných pozicích, a tím zredukovat paměťovou a výpočetní náročnost
- technologii ACELP využívá cca
 - 2,4 miliard uživatelů mobilních telefonů na celém světě
 - 35 milionů uživatelů přehrávačů MP3
 - 500 milionů uživatelů internetových přehrávačů RealPlayer nebo MediaPlayer

Prohledávání fixní knihovny

1 podrámec řeči (~5ms)

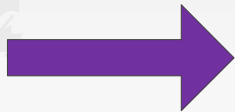


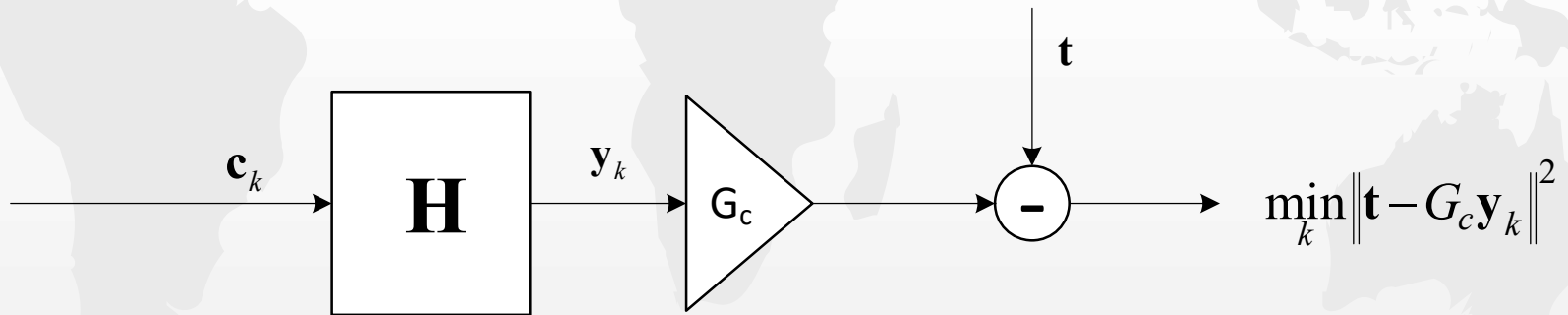
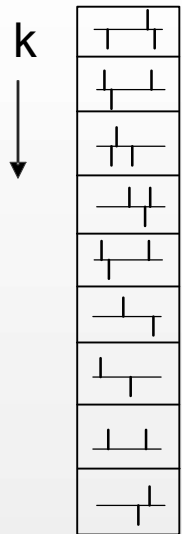
Zavedení algebraické knihovny



algebraická knihovna (až 80 bitů)

Prohledávání fixní knihovny

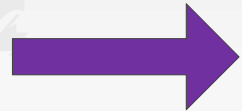
$$\min_k \|\mathbf{t} - G_c \mathbf{y}_k\|^2 \quad \frac{\partial}{\partial G_c}$$




algebraická knihovna (až 80 bitů)

Prohledávání algebraické knihovny

$$\min_k \|\mathbf{t} - G_c \mathbf{y}_k\|^2$$

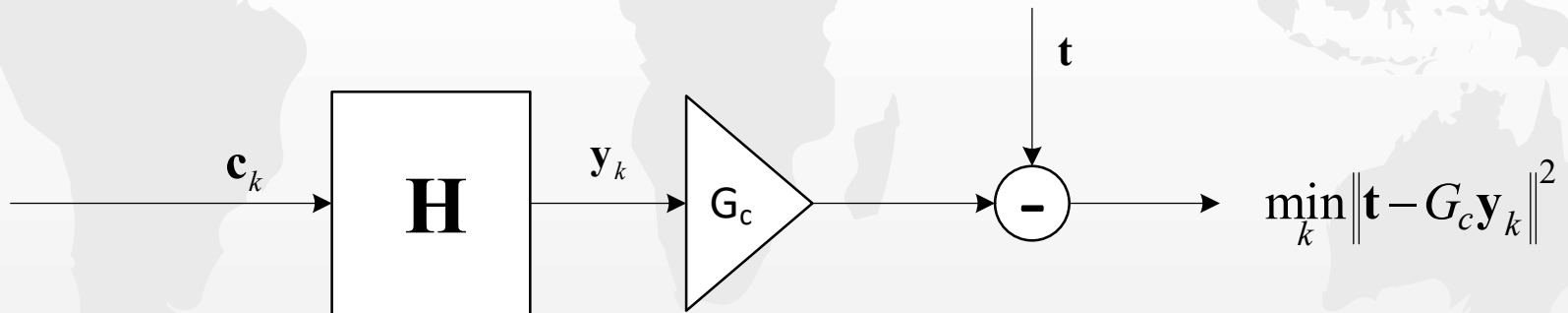
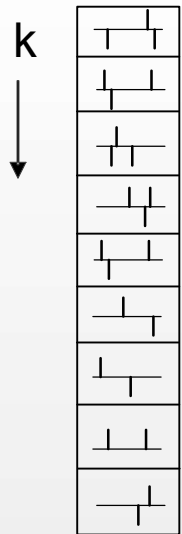


$$\frac{\partial}{\partial G_c}$$

$$\max_k \frac{\mathbf{t}^T \cdot \mathbf{y}_k}{\mathbf{y}_k^T \cdot \mathbf{y}_k}$$

korelace mezi cílovým (target) vektorem a testovaným vektorem

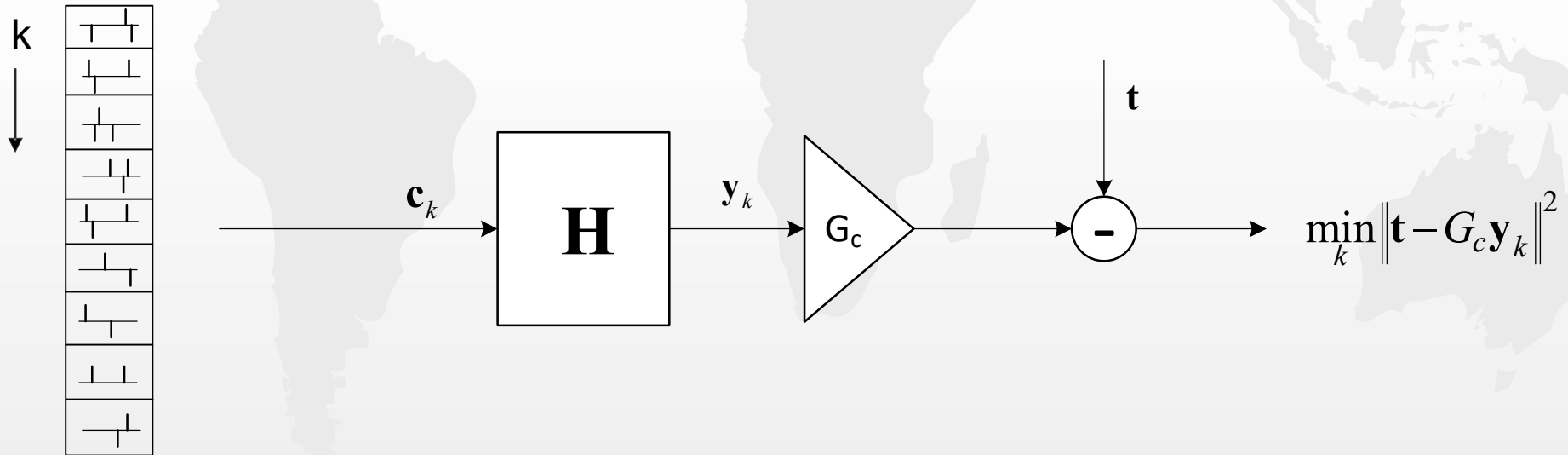
energie testovaného vektoru



algebraická knihovna (až 80 bitů)

Prohledávání algebraické knihovny

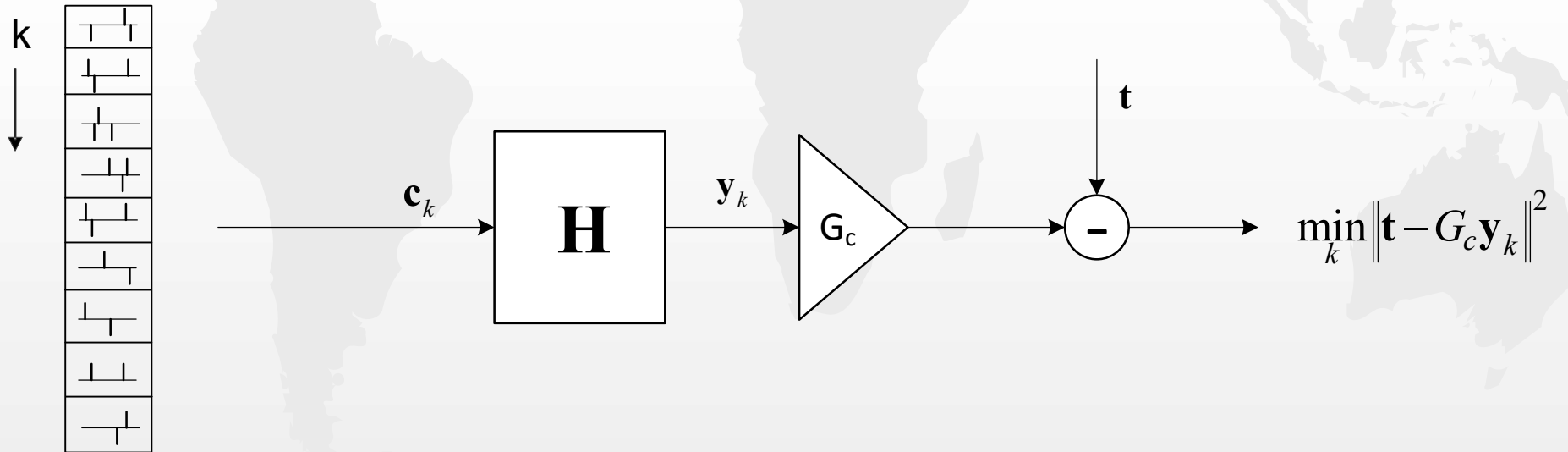
$$\min_k \|\mathbf{t} - G_c \mathbf{y}_k\|^2 \xrightarrow{\frac{\partial}{\partial G_c}} \max_k \frac{(\mathbf{t}^T \cdot \mathbf{y}_k)^2}{\mathbf{y}_k^T \cdot \mathbf{y}_k} \xrightarrow{\quad} \max_k \frac{(\mathbf{t}^T \cdot \mathbf{H} \cdot \mathbf{c}_k)^2}{\mathbf{c}_k^T \cdot \mathbf{H}^T \cdot \mathbf{H} \cdot \mathbf{c}_k}$$



algebraická knihovna (až 80 bitů)

Prohledávání algebraické knihovny

$$\min_k \|\mathbf{t} - G_c \mathbf{y}_k\|^2 \xrightarrow{\frac{\partial}{\partial G_c}} \max_k \frac{(\mathbf{t}^T \cdot \mathbf{y}_k)^2}{\mathbf{y}_k^T \cdot \mathbf{y}_k} \xrightarrow{\quad} \max_k \frac{(\mathbf{t}^T \cdot \mathbf{H} \cdot \mathbf{c}_k)^2}{\mathbf{c}_k^T \cdot \mathbf{H}^T \cdot \mathbf{H} \cdot \mathbf{c}_k} \xrightarrow{\quad} \max_k \frac{(\mathbf{d}^T \cdot \mathbf{c}_k)^2}{\mathbf{c}_k^T \cdot \mathbf{\Phi} \cdot \mathbf{c}_k}$$



algebraická knihovna (až 80 bitů)

Prohledávání algebraické knihovny

$$\max_k \frac{(\mathbf{d}^T \cdot \mathbf{c}_k)^2}{\mathbf{c}_k^T \cdot \Phi \cdot \mathbf{c}_k}$$

Lze prohledávat rychle, pokud \mathbf{c}_k obsahuje jen velmi málo nenulových prvků s hodnotami +1 nebo -1

$$\mathbf{d}^T \cdot \mathbf{c}_k$$

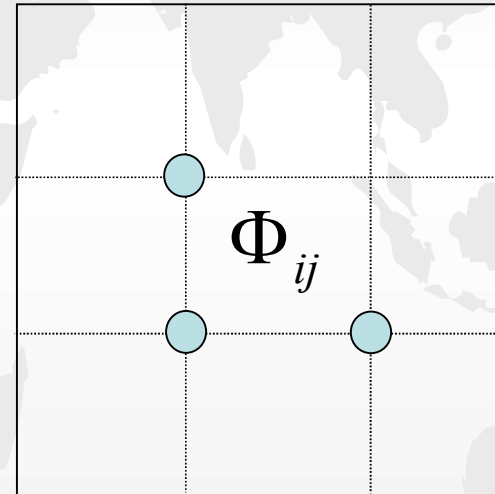
$$d_0 \ d_1 \ d_2 \ \dots \ d_9$$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$= d_3 + d_6$$

$$\mathbf{c}_k^T \cdot \Phi \cdot \mathbf{c}_k$$

$$0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0$$



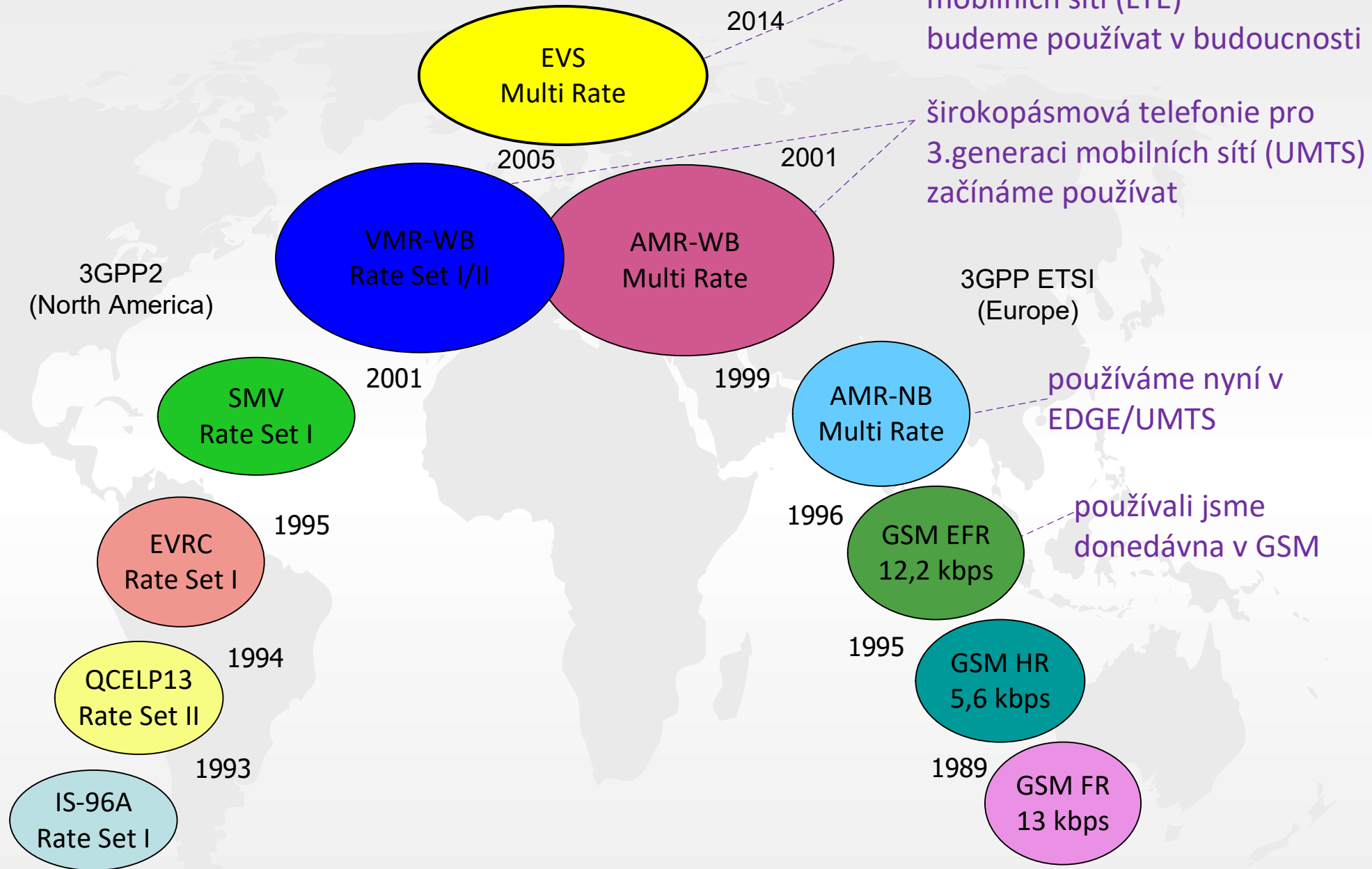
$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$= \Phi_{3,3} + \Phi_{6,6} + 2\Phi_{3,6}$$

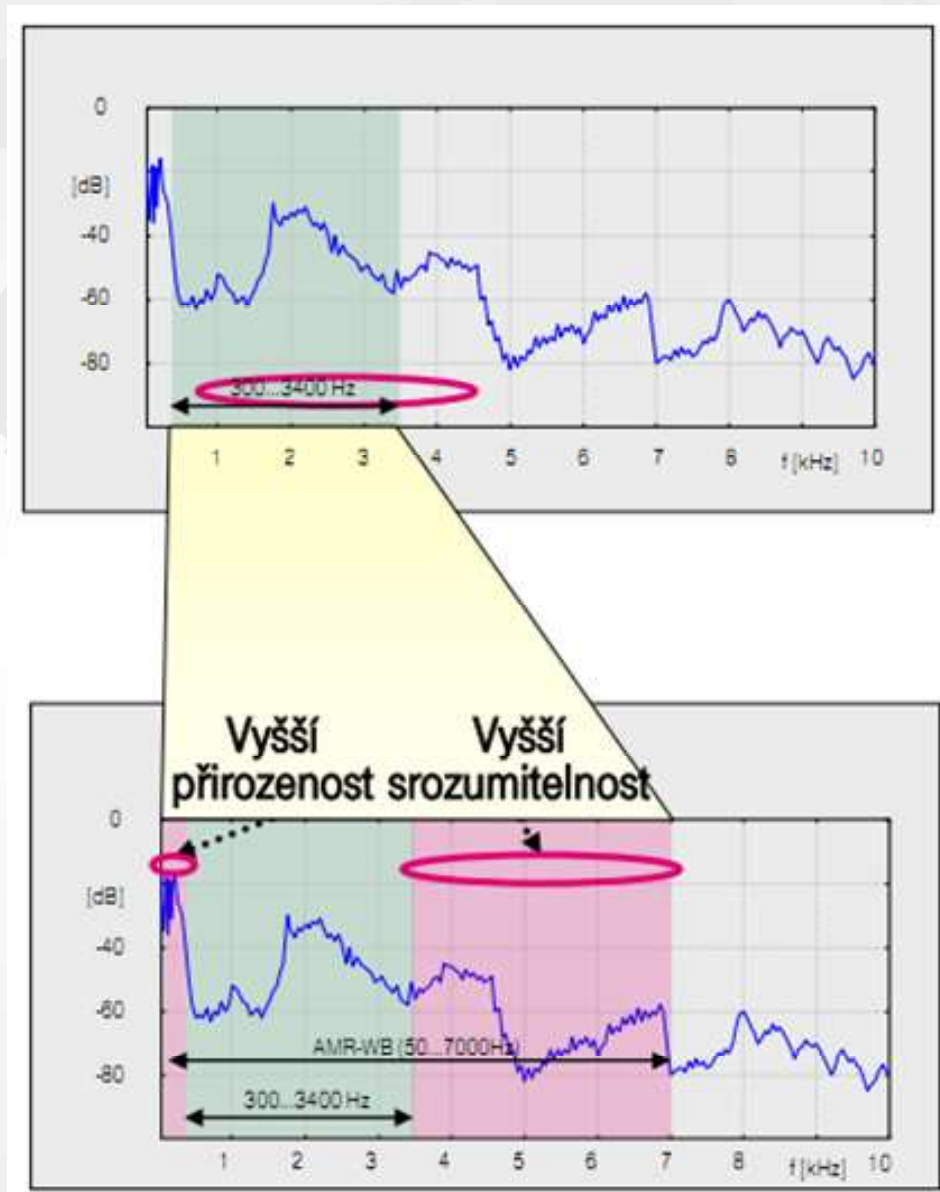


ACELP ve světě

Technologie ACELP v mezinárodních standardech



Od AMR-NB k AMR-WB (HD VOICE)



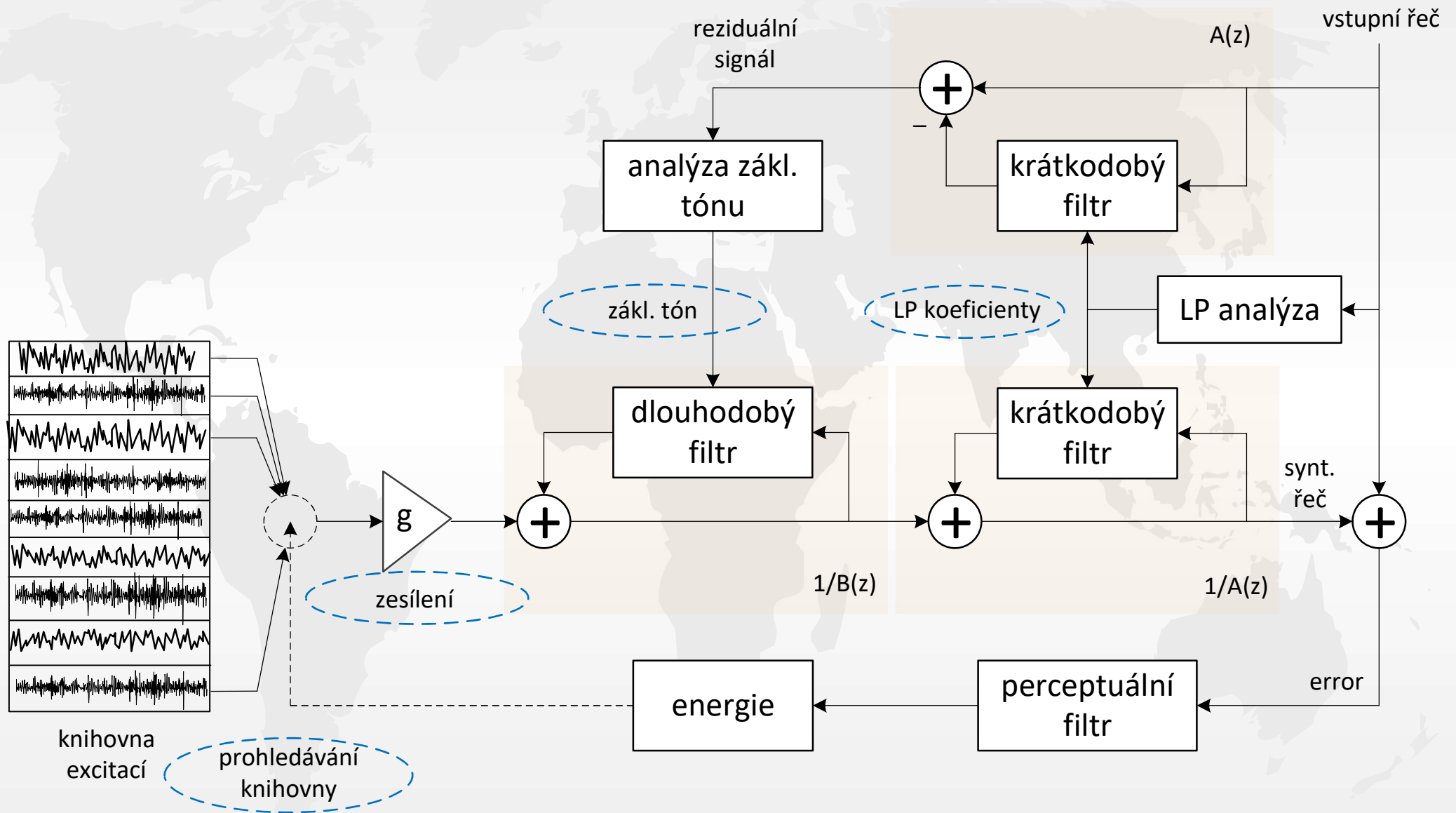
- HD voice demo na <https://www.youtube.com/watch?v=Y4bb3b9PiRg>

HD VOICE



Konec

Základní schéma kodeku



Základní schéma kodeku

