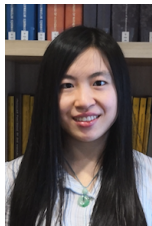


Fast Beam Alignment via Pure Exploration in Multi-Armed Bandits

Yi Wei



Zixin Zhong

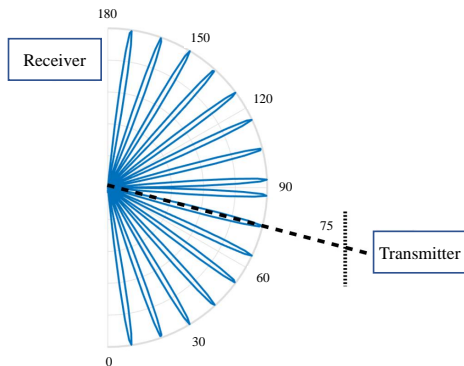


Vincent Y. F. Tan

National University of Singapore

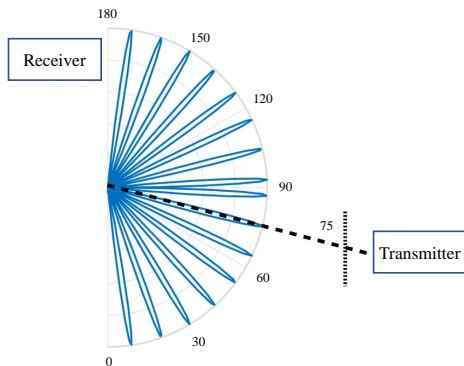
June 29, 2022

The Beam Alignment Problem



- mmWave Communications
- Beams at Tx and Rx are **narrow directional**

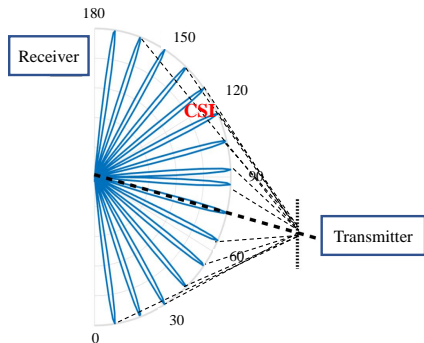
The Beam Alignment Problem



- mmWave Communications
- Beams at Tx and Rx are **narrow directional**

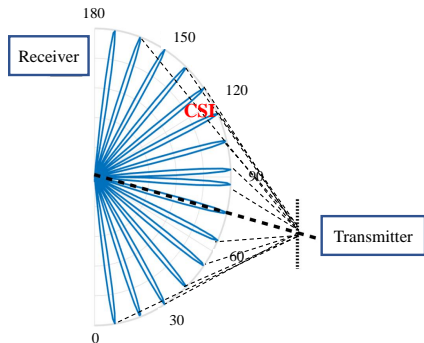
Beam alignment (BA) is to ensure the **transmitter** and **receiver** beams are **accurately aligned** to establish a reliable communication link.

Some Fundamental Challenges in Beam Alignment



- To find optimal Rx/Tx pair, **CSI** corresponding to each pair is measured
- **Frequency** of each measurement is **high**
- **High** beam alignment **latency**

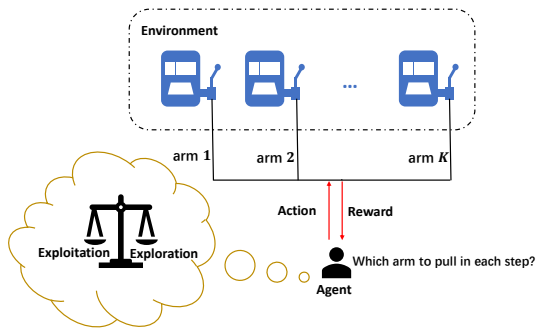
Some Fundamental Challenges in Beam Alignment



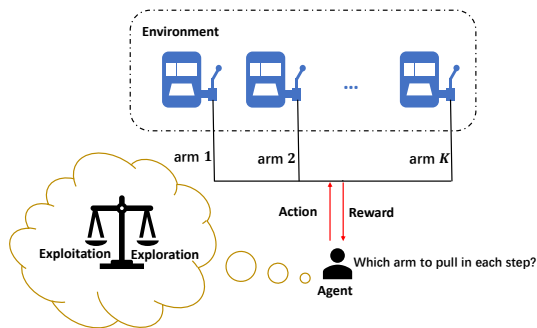
- To find optimal Rx/Tx pair, **CSI** corresponding to each pair is measured
- **Frequency** of each measurement is **high**
- **High** beam alignment **latency**

Beam alignment latency will **increase** with the **number of antennas** at the receivers and transmitters.

Multi-Armed Bandits

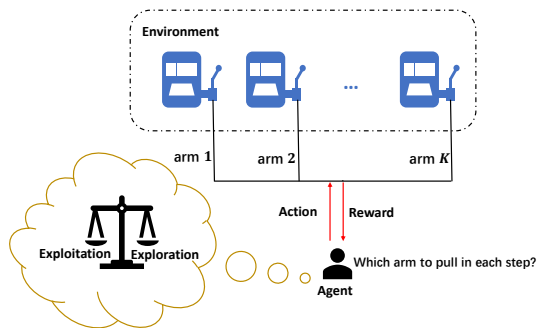


Multi-Armed Bandits



Pure Exploration: Identify the best arm (arm with largest mean) as quickly as possible.

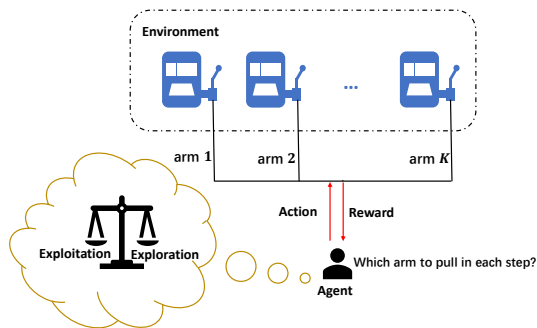
Multi-Armed Bandits



Pure Exploration: Identify the best arm (arm with largest mean) as quickly as possible.

Main Contribution: Formulate and solve the beam alignment problem using ideas in pure exploration in the fixed-confidence setting.

Multi-Armed Bandits



Pure Exploration: Identify the best arm (arm with largest mean) as quickly as possible.

Main Contribution: Formulate and solve the beam alignment problem using ideas in pure exploration in the fixed-confidence setting.

Exploit **structure** in the beam alignment problem.

System Model

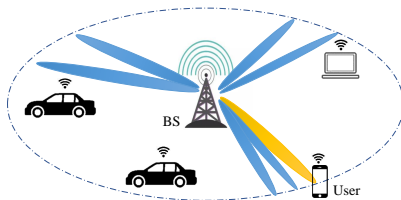


Figure: A mmWave massive MISO system system.

- **Massive mmWave MISO system:** A base station (BS) equipped with N transmit antennas serves a single-antenna user

System Model

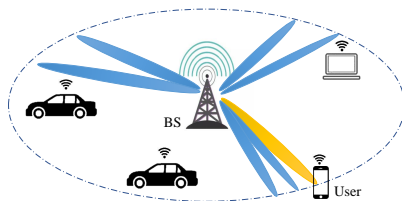


Figure: A mmWave massive MISO system.

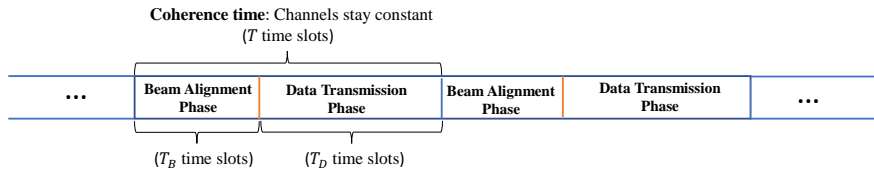
- **Massive mmWave MISO system:** A base station (BS) equipped with N transmit antennas serves a single-antenna user
- Adopt the widely-used Saleh–Valenzuela channel model

$$\mathbf{h} = \beta^{(1)} \mathbf{a}(\theta^{(1)}) + \sum_{l=2}^L \beta^{(l)} \mathbf{a}(\theta^{(l)})$$

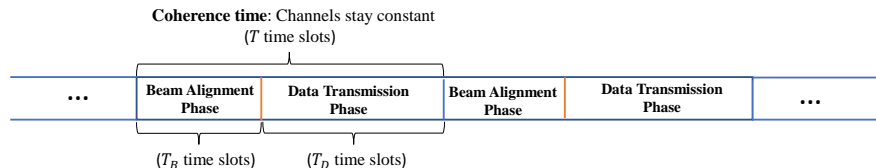
1 line-of-sight (LoS) path \geq L - 1 non-LoS (NLoS) paths

Amplitude

Transmission Scheme



Transmission Scheme



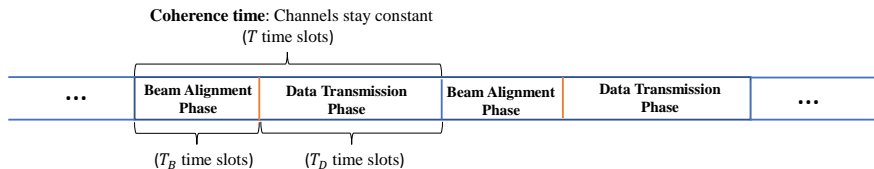
- **Beam alignment phase:** Fast beam alignment algorithm searches the optimal beam from a given codebook

$$\mathcal{C}, \{ \mathbf{f}_k = \mathbf{a}(-1 + 2k/K) \mid k = 0, 1, \dots, K - 1 \}$$

where the array response vector:

$$\mathbf{a}(x) = \frac{1}{\sqrt{N}} [1, e^{j^2-dx}, e^{j^2-2dx}, \dots, e^{j^2-(N-1)dx}]^H \in \mathbb{C}^{N \times 1}$$

Transmission Scheme



- Beam alignment phase:** Fast beam alignment algorithm searches the optimal beam from a given codebook

$$\mathcal{C}, \{ \mathbf{f}_k = \mathbf{a}(-1 + 2k/K) \mid k = 0, 1, \dots, K - 1 \}$$

where the array response vector:

$$\mathbf{a}(x) = \frac{1}{\sqrt{N}} [1, e^{j^2-dx}, e^{j^2-2dx}, \dots, e^{j^2-(N-1)dx}]^H \in \mathbb{C}^{N \times 1}$$

- Data transmission phase:** BS transmits effective data using the selected \mathbf{f} . Received signal at the user in time slot t is

$$y_t = \bar{p} \mathbf{h}^H \mathbf{f} s_t + n_t,$$

Beam Alignment Phase

- **System Throughput Performance:** Effective achievable rate

$$R_{\text{eff}} = \left(1 - \frac{T_B}{T_D}\right) \log \left(1 + \frac{p/h^H f \beta^2}{2}\right)$$

Time to search for optimal beam T_B to be minimized for high R_{eff}

Beam Alignment Phase

- **System Throughput Performance:** Effective achievable rate

$$R_{\text{eff}} = 1 - \frac{T_B}{T_D} \log \left(1 + \frac{p|\mathbf{h}^H \mathbf{f}|^2}{2} \right)$$

Time to search for optimal beam T_B to be minimized for high R_{eff}

- **Measurement:** The received signal power

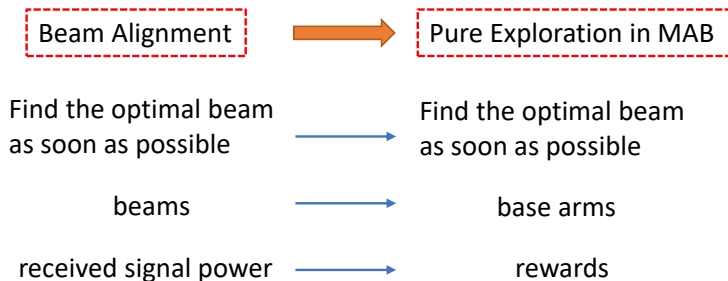
$$R(\mathbf{f}_k) = |\sqrt{p}\mathbf{h}^H \mathbf{f}_k + n|^2 = p|\mathbf{h}^H \mathbf{f}_k|^2 + 2\sqrt{p}\Re(\mathbf{h}^H \mathbf{f}_k n^*) + |n|^2$$

Heteroscedastic Gaussian Variable $\mathcal{N}(p|\mathbf{h}^H \mathbf{f}_k|^2, 2p|\mathbf{h}^H \mathbf{f}_k|^2\sigma^2)$ **Gamma Variable** $\Gamma(1, 1/\sigma^2)$

Approximate (Because: noise power \ll transmit power)

$$r_k = p|\mathbf{h}^H \mathbf{f}_k|^2 + 2\sqrt{p}\Re(\mathbf{h}^H \mathbf{f}_k n^*)$$

Relation to MABs and Properties



Relation to MABs and Properties

Beam Alignment



Pure Exploration in MAB

Find the optimal beam
as soon as possible



Find the optimal beam
as soon as possible

beams



base arms

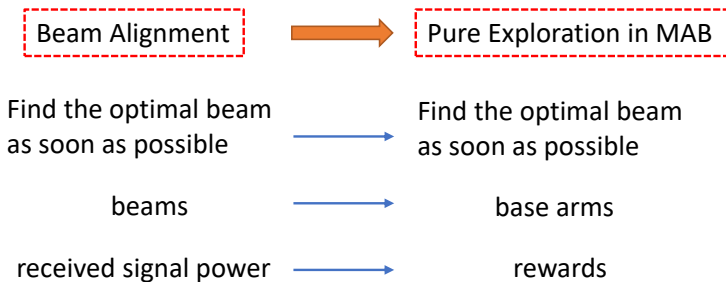
received signal power



rewards

Property: Let $\mu = (\mu_1, \dots, \mu_K)$, and let $\mu_{(1)} \leq \mu_{(2)} \leq \mu_{(3)} \leq \dots \leq \mu_{(K)}$ be the sorted sequence of means.

Relation to MABs and Properties



Property: Let $\mu = (\mu_1, \dots, \mu_K)$, and let $\mu_{(1)} \leq \mu_{(2)} \leq \mu_{(3)} \leq \dots \leq \mu_{(K)}$ be the sorted sequence of means.

1. The **means of the rewards** associated with **close-by** arms are **close**.
2. The **variance** of the reward of an arm is linearly related to its **mean**

$$\sigma_k^2 = 2\mu_k^2.$$

Beam Codebooks Possess the Group Property

$\frac{1}{J}$ -lower resolution beam codebook

- Constructed by grouping the nearby beams in the codebook \mathcal{C}

$$\mathcal{C}_{(J)}, \quad \mathbf{b}_g = \left[\mathbf{f}_k \right]_{k=J(g-1)+1}^{Jg} : g = 0, 1, \dots, G-1$$

Beam Codebooks Possess the Group Property

$\frac{1}{J}$ -lower resolution beam codebook

- Constructed by grouping the nearby beams in the codebook \mathcal{C}

$$\mathcal{C}_{(J)}, \quad \mathbf{b}_g = \sum_{k=J(g-1)+1}^{Jg} \mathbf{f}_k : g = 0, 1, \dots, G-1$$

- Received power for beam \mathbf{b}_g

$$R_g = p / \mathbf{h}^H \mathbf{b}_g \mathbf{b}_g^H + 2 \bar{p} (\mathbf{h}^H \mathbf{b}_g n),$$

follows a **heteroscedastic Gaussian distribution**.

Beam Codebooks Possess the Group Property

$\frac{1}{J}$ -lower resolution beam codebook

- Constructed by grouping the nearby beams in the codebook \mathcal{C}

$$\mathcal{C}_{(J)}, \quad \mathbf{b}_g = \left\{ \mathbf{f}_k : k = J(g-1)+1, \dots, Jg \right\}, \quad g = 0, 1, \dots, G-1$$

- Received power for beam \mathbf{b}_g

$$R_g = p / \mathbf{h}^H \mathbf{b}_g \mathbf{b}_g^H + 2 \bar{p} (\mathbf{h}^H \mathbf{b}_g n),$$

follows a **heteroscedastic Gaussian distribution**.

- Information of a **set of beams** can be obtained at each time step

Problem Setup for Bandit Beam Alignment

Bandit BA Problem

- K **base arms** $[K]$, $\{1, \dots, K\}$: associated with the beam \mathbf{f}_k
- $\{[K], J\}$: set of all non-empty consecutive tuples of length J
 $\{[6], 2\} = \{\{1\}, \{1, 2\}, \{2\}, \{2, 3\}, \{3\}, \{3, 4\}, \{4\}, \{4, 5\}, \{5\}, \{5, 6\}, \{6\}\}$
- (K, J) -**super arm**: each tuple in $\{[K], J\}$, associated with a “grouped beam” $\mathbf{b}_g \in \mathcal{C}_{(J)}$.

Bandit BA Problem Setup

In time step t

- Choose an action (or a (K, J) -super arm) $A(t) \in \{[K], J\}$

Bandit BA Problem Setup

In time step t

- Choose an action (or a (K, J) -super arm) $A(t) \in \{[K], J\}$
- Observe the reward

$$R(t) = F(\mathbf{f}_k, \rho, \mathbf{h}, n_t)$$

$k = A(t)$

where $F(\mathbf{f}, \rho, \mathbf{h}, n) = \rho \|\mathbf{h}^H \mathbf{f}\|^2 + 2 \bar{\rho} (\mathbf{h}^H \mathbf{f} n)$

Bandit BA Problem Setup

In time step t

- Choose an action (or a (K, J) -super arm) $A(t) \in \{[K], J\}$
- Observe the reward

$$R(t) = F_{k, A(t)}(\mathbf{f}_k, \rho, \mathbf{h}, n_t)$$

where $F(\mathbf{f}, \rho, \mathbf{h}, n) = \rho / \mathbf{h}^H \mathbf{f} \rho + 2 \bar{\rho} (\mathbf{h}^H \mathbf{f} n)$

- Each observed reward $R(t)$ also follows a **heteroscedastic Gaussian distribution**.

Bandit BA Problem Setup

Algorithm: $\mathcal{A} := \{(A(t))_t, \dots, J\}$

- **Sampling rule** \mathcal{A}_t : Determines the $([K], J)$ -super arm $A(t)$ to pull at time step t based on the observation history and the arm history $\{A(1), R(1), A(2), R(2), \dots, A(t-1), R(t-1)\}$
- **Stopping rule**: Leads to a stopping time τ satisfying $P(\tau < \infty) = 1$
- **Recommendation rule** \mathcal{A}_τ : Outputs a base arm $k \in [K]$.

Bandit BA Problem Setup

Algorithm: $\pi := \{(\pi_t)_t, \dots, J\}$

- **Sampling rule** π_t : Determines the $([K], J)$ -super arm $A(t)$ to pull at time step t based on the observation history and the arm history $\{A(1), R(1), A(2), R(2), \dots, A(t-1), R(t-1)\}$
- **Stopping rule**: Leads to a stopping time τ satisfying $P(\tau < \infty) = 1$
- **Recommendation rule** \hat{k} : Outputs a base arm $k \in [K]$.

Aim: Use as few samples as possible to output a guess of the optimal arm with probability at least $1 - \epsilon$.

Some Notations for the General Lower Bound

- Heteroscedastic Gaussian bandit instance:

$$= N(\mu_1, 2\mu_1^2), \dots, N(\mu_K, 2\mu_K^2)$$

- Optimal arm A^* :

$$\mu_{A^*} = \operatorname{argmax}_{k \in [K]} \mu_k$$

- Set of probability distributions

$$W_K := \{ \mathbf{w} \in \mathbb{R}_+^K : \sum_{k=1}^K w_k = 1 \}$$

- Alternative Set

$$\text{Alt}(A^*) := \{ \mathbf{u} \in W_K : A(\mathbf{u}) = A^* \}$$

General Lower Bound

Theorem

For any $(\mathcal{U}, \mathcal{J})$ -PAC algorithm where $\mathcal{U} \subseteq (0, 1)$, it holds that

$$E[\text{error}] \geq c(\mathcal{U}) \ln \frac{1}{4},$$

where

$$c(\mathcal{U})^{-1} = \sup_{\mathbf{w}} \inf_{\text{Alt}(\mathcal{U})} \sum_{k=1}^K w_k D_{\text{HG}}(\mu_k, \mu_k^{\mathbf{u}})$$

and the *heteroscedastic relative entropy* is

$$D_{\text{HG}}(\mu_i, \mu_j) = \frac{1}{2} \ln \frac{\mu_j}{\mu_i} + \frac{\mu_i}{2\mu_j} + \frac{(\mu_j - \mu_i)^2}{4\mu_j^2} - \frac{1}{2}.$$

Two-Phase Heteroscedastic Track & Stop (2PHT&S)

Main Idea

to exploit the **prior knowledge** which have not been considered by existing bandit-based beam alignment algorithms:

- Close-by correlation
- Heteroscedasticity
- Group property

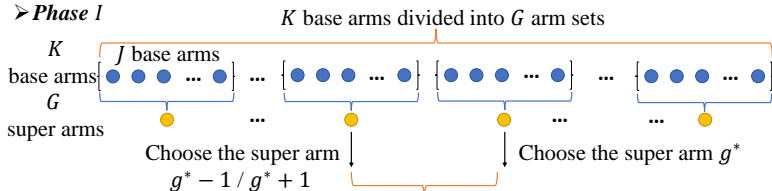
Two-Phase Heteroscedastic Track & Stop (2PHT&S)

Main Idea

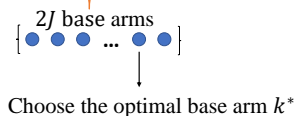
to exploit the **prior knowledge** which have not been considered by existing bandit-based beam alignment algorithms:

- Close-by correlation
- Heteroscedasticity
- Group property

➤ Phase I



➤ Phase II



2PHT&S Algorithm

Phase I: Search for the optimal super arm w.p. 1 1

2PHT&S Algorithm

Phase I: Search for the optimal super arm w.p. $1 - \epsilon$

Group K base arms into G arm sets to reduce the search space

2PHT&S Algorithm

Phase I: **Search for the optimal super arm w.p. $1 - \epsilon$**

Group K base arms into G arm sets to **reduce the search space**

At each time,

Choose one super arm (beam group) by the sampling rule in a new **Heteroscedastic Track and Stop (HT&S) algorithm**

2PHT&S Algorithm

Phase I: Search for the optimal super arm w.p. $1 - \epsilon$

Group K base arms into G arm sets to reduce the search space

At each time,

Choose one super arm (beam group) by the sampling rule in a new Heteroscedastic Track and Stop (HT&S) algorithm

Use the grouped beam to transmit the pilot symbols and observe the grouped reward $R_g(t)$.

2PHT&S Algorithm

Phase I: Search for the optimal super arm w.p. $1 - \epsilon$

Group K base arms into G arm sets to reduce the search space

At each time,

Choose one super arm (beam group) by the sampling rule in a new Heteroscedastic Track and Stop (HT&S) algorithm

Use the grouped beam to transmit the pilot symbols and observe the grouped reward $R_g(t)$.

Select the optimal super arm

$$g = \underset{g \in [G]}{\operatorname{argmax}} E[R_g(t)]:$$

Phase II: Search for the optimal base arm w.p. 1 2

Phase II: Search for the optimal base arm w.p. $\frac{1}{2}$

Construct a base arm set (including the optimal super arm and its neighboring super arm)

Phase II: Search for the optimal base arm w.p. $\frac{1}{2}$

Construct a base arm set (including the optimal super arm and its neighboring super arm)

Search for the optimal base arm in the base arm set using the HT&S algorithm

HT&S: An Improved Track & Stop Algorithm

Sampling Rule: Estimate the number of times each arm should be sampled

$$Q(t) = \begin{cases} \operatorname{argmin}_i T_i(t-1); & \text{if } \min_i T_i(t-1) \leq \bar{p}_t; \\ \operatorname{argmax}_i t \hat{w}_i(t-1) T_i(t-1); & \text{otherwise:} \end{cases}$$

HT&S: An Improved Track & Stop Algorithm

Sampling Rule: Estimate the number of times each arm should be sampled

$$Q(t) = \begin{cases} \operatorname{argmin}_i T_i(t-1); & \text{if } \min_i T_i(t-1) \leq p_{\bar{t}}; \\ \operatorname{argmax}_i t \hat{w}_i(t-1) T_i(t-1); & \text{otherwise:} \end{cases}$$

Stopping Rule: Stop when the number of pulls of all arms meet a certain requirement

Threshold to be $(t; \delta) = \ln(t/\delta)$, and the stopping rule is

$$= \min_{t \geq 2N} Z(t) \leq (t; \delta) :$$

HT&S: An Improved Track & Stop Algorithm

- **Sampling Rule:** Estimate the number of times each arm should be sampled

$$Q(t) = \begin{cases} \operatorname{argmin}_i T_i(t-1), & \text{if } \min_i T_i(t-1) \leq \bar{t}, \\ \operatorname{argmax}_i t \hat{w}_i(t-1) - T_i(t-1), & \text{otherwise.} \end{cases}$$

- **Stopping Rule:** Stop when the number of pulls of all arms meet a certain requirement

Threshold to be $(t, \epsilon) = \ln(t/\epsilon)$, and the stopping rule is

$$= \min_{t \in \mathbb{N}} : Z(t) \leq (t, \epsilon).$$

- **Heteroscedasticity:** Considered explicitly in the calculations of estimated reward $\hat{w}_i(t-1)$ and statistic $Z(t)$

Time Complexity Analysis of 2PHT&S

Theorem

Let the means of the *super* and *base arms* be respectively

$$\begin{aligned} \mathbf{s} &:= N(\mu_1^s, 2\mu_1^s{}^2), \dots, N(\mu_G^s, 2\mu_G^s{}^2) \quad \text{and} \\ \mathbf{b} &:= N(\mu_{S_f(1)}^b, 2\mu_{S_f(1)}^b{}^2), \dots, N(\mu_{S_f(2J)}^b, 2\mu_{S_f(2J)}^b{}^2) \end{aligned}$$

where

$$\mu_g^s = p \mathbf{h}^H \mathbf{f}_k \quad k \in S_g$$

Time Complexity Analysis of 2PHT&S

Theorem

Let the means of the *super* and *base arms* be respectively

$$\begin{aligned} \mathbf{s} &:= N(\mu_1^s, 2\mu_1^s{}^2), \dots, N(\mu_G^s, 2\mu_G^s{}^2) \quad \text{and} \\ \mathbf{b} &:= N(\mu_{S_f(1)}^b, 2\mu_{S_f(1)}^b{}^2), \dots, N(\mu_{S_f(2J)}^b, 2\mu_{S_f(2J)}^b{}^2) \end{aligned}$$

where

$$\mu_g^s = p \mathbf{h}^H \mathbf{f}_k \quad k \in S_g$$

Under 2PHT&S, we have

$$\limsup_0 \frac{\mathbb{E}[2PHT\&S]}{\ln(1/\epsilon)} = C_s + C_b,$$

where C_s and C_b represent time complexities of Phases I and II resp.

Simulation Results

Experiment Setup

Massive mmWave MISO system

Base station equipped with 64 transmit antennas serving a single-antenna user

Size of codebook is set as $K = 128$.

Simulation Results

Experiment Setup

Massive mmWave MISO system

Base station equipped with 64 transmit antennas serving a single-antenna user

Size of codebook is set as $K = 128$.

Baseline Algorithms

Original Track-and-Stop (T&S) algorithm (Garivier et al. 2016)

Two-Phase Track-and-Stop (2PT&S) algorithm

Heteroscedastic Track-and-Stop (HT&S) algorithm

A. Garivier and E. Kaufmann, "Optimal best arm identification with fixed confidence," in PMLR, 2016, pp. 998-1027.

Numerical Simulations on Synthetic Data

Figure: Mean of the reward of each base arm and super arm ($p = 10$ dBm).

Table: Average time complexities over 100 experiments when $\epsilon = 0:1$

Power	4		6		8		10		12	
T&S	1154.3	338:7	654.6	212:1	382.5	129:6	209.4	68:6	133.7	8:9
HT&S	473.2	275:5	271.4	143:4	175.6	69:2	133.2	24:1	123.9	6:5
2PT&S	206.2	60:4	120.2	35:0	68.4	19:4	49.1	4:6	45.2	1:1
2HPT&S	84.3	41:5	58.0	19:6	48.4	6:3	45.5	1:6	45	0

Numerical Simulations on Synthetic Data

Figure: Mean of the reward of each base arm and super arm ($p = 10$ dBm).

Table: Average time complexities over 100 experiments when $\epsilon = 0:1$

Power	4		6		8		10		12	
T&S	1154.3	338:7	654.6	212:1	382.5	129:6	209.4	68:6	133.7	8:9
HT&S	473.2	275:5	271.4	143:4	175.6	69:2	133.2	24:1	123.9	6:5
2PT&S	206.2	60:4	120.2	35:0	68.4	19:4	49.1	4:6	45.2	1:1
2HPT&S	84.3	41:5	58.0	19:6	48.4	6:3	45.5	1:6	45	0

Experiments on [real data for a practical scenario in a city](#) available in the longer version of our paper.

Conclusions

- Formulate the **beam alignment problem** as a multi-armed bandit problem under the **pure exploration** setting

Conclusions

- Formulate the **beam alignment problem** as a multi-armed bandit problem under the **pure exploration** setting
- Derived a general lower bound on the sample complexity

Conclusions

- Formulate the **beam alignment problem** as a multi-armed bandit problem under the **pure exploration** setting
- Derived a general lower bound on the sample complexity
- Exploited the **structure** and **properties** of the beam alignment problem to derive an algorithm **2PHT&S** with a reduced time complexity

Conclusions

- Formulate the **beam alignment problem** as a multi-armed bandit problem under the **pure exploration** setting
- Derived a general lower bound on the sample complexity
- Exploited the **structure** and **properties** of the beam alignment problem to derive an algorithm **2PHT&S** with a reduced time complexity
- Simulations demonstrate **superior performances** over benchmarks